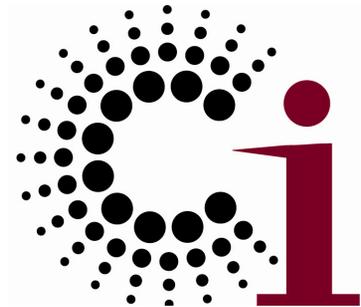


Why Computer Science is Fundamental to Almost Everything

Ian Foster



Computation Institute

Argonne National Lab & University of Chicago

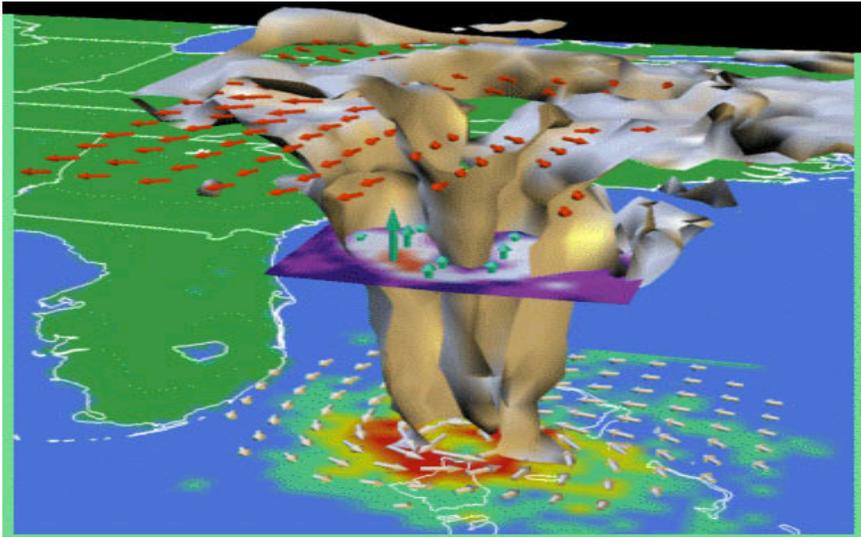


“Applied computer science is now playing the role that mathematics did from the 17th through the 20th centuries: providing an orderly, formal framework & exploratory apparatus for other sciences.”

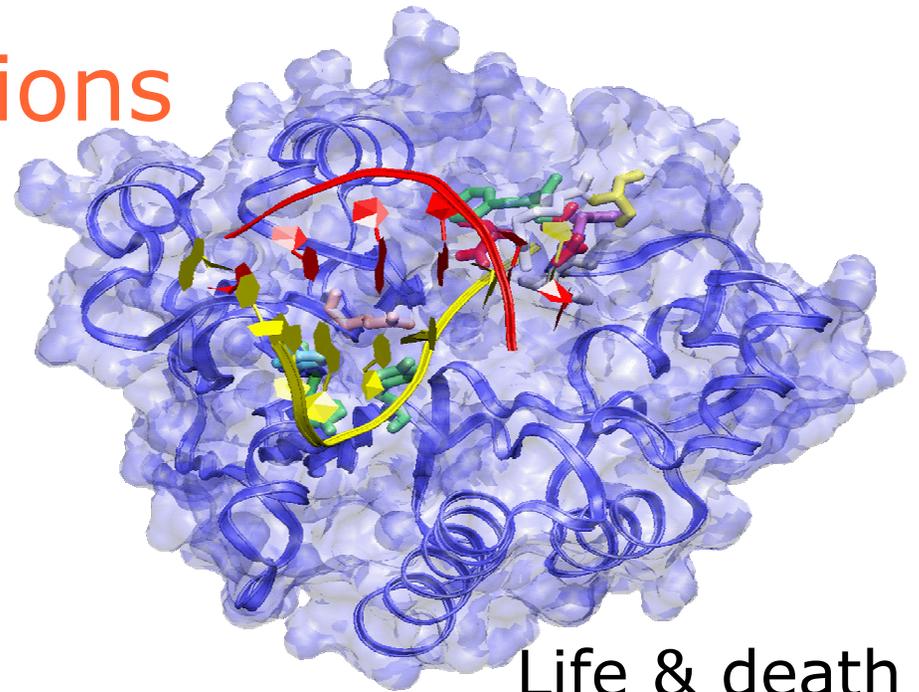
—George Djorgovski



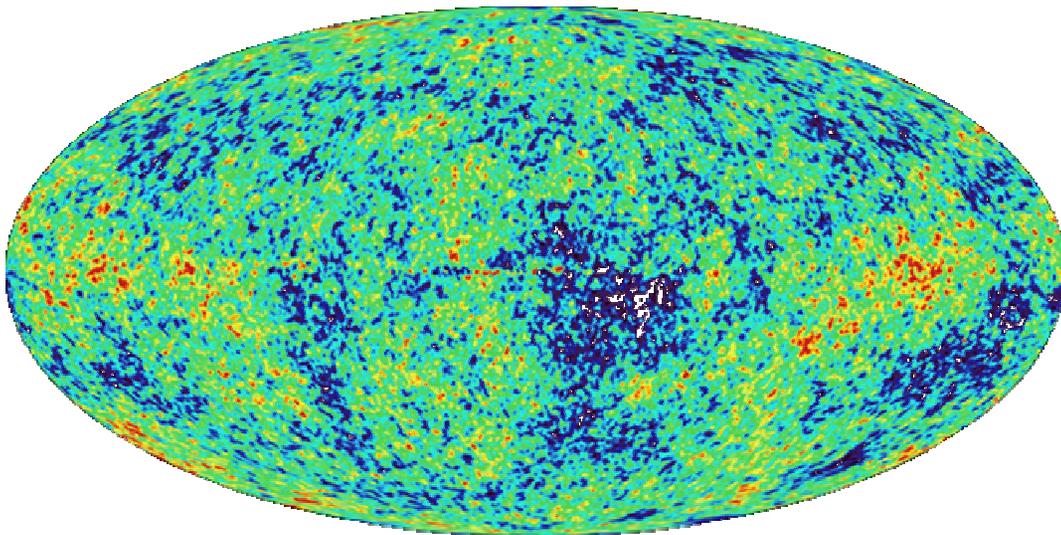
The Big Questions



Future of the planet



Life & death



Nature of the universe



Consciousness

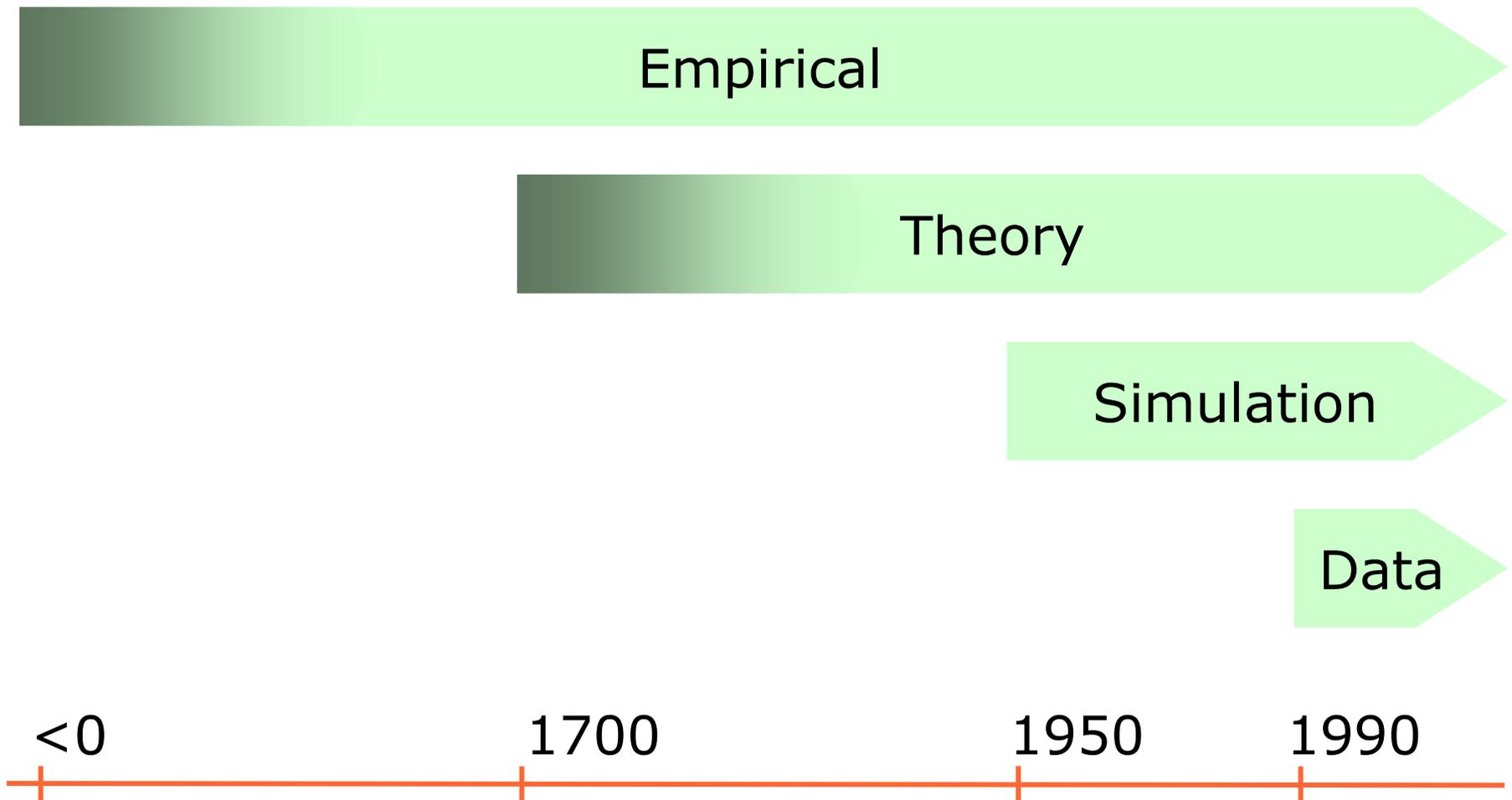


The Little Questions

- Friends
- Sales
- Entertainment
- Spelling
- Parking



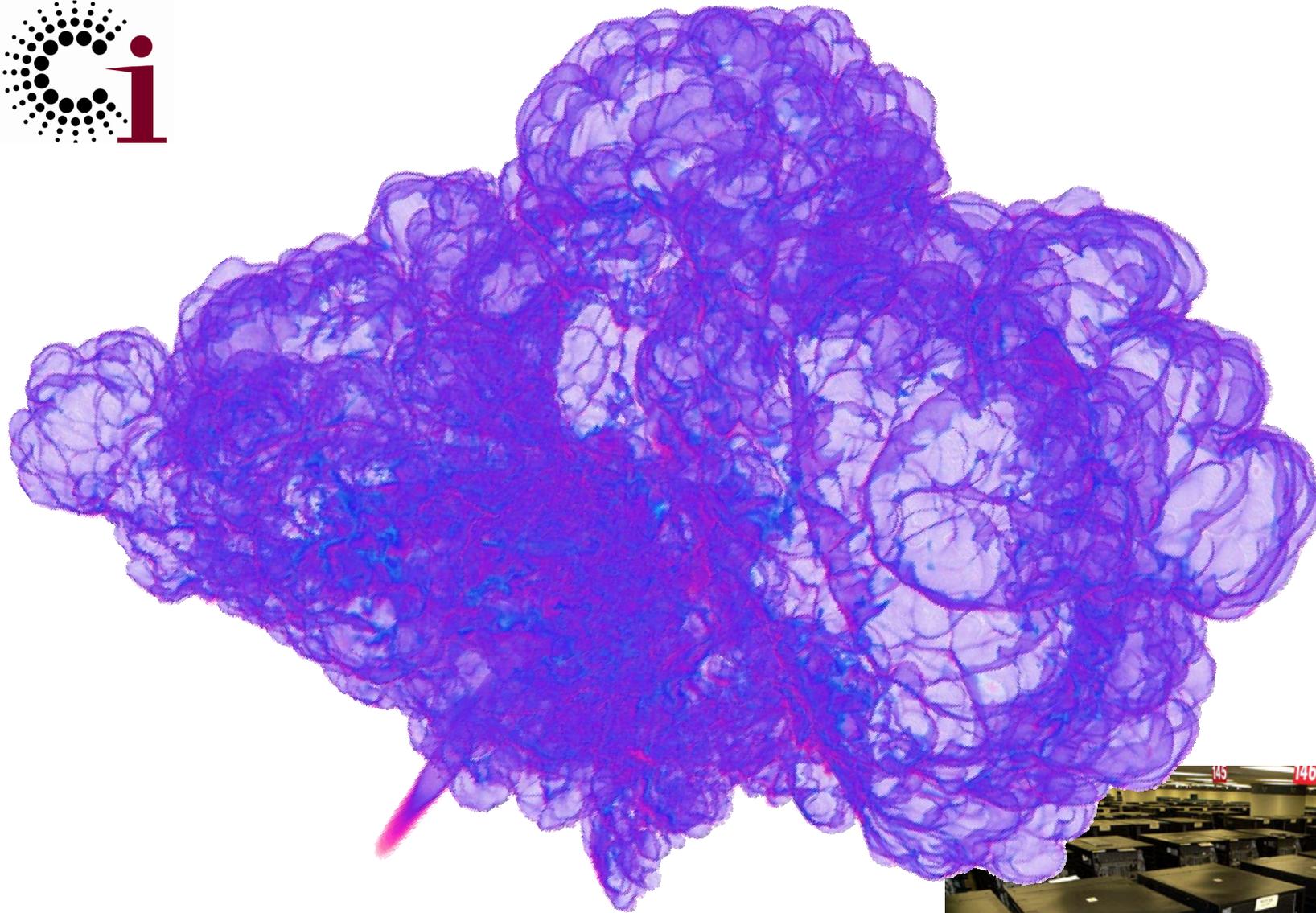
How Do We Answer Them?





Type Ia Supernova: SN 1994D



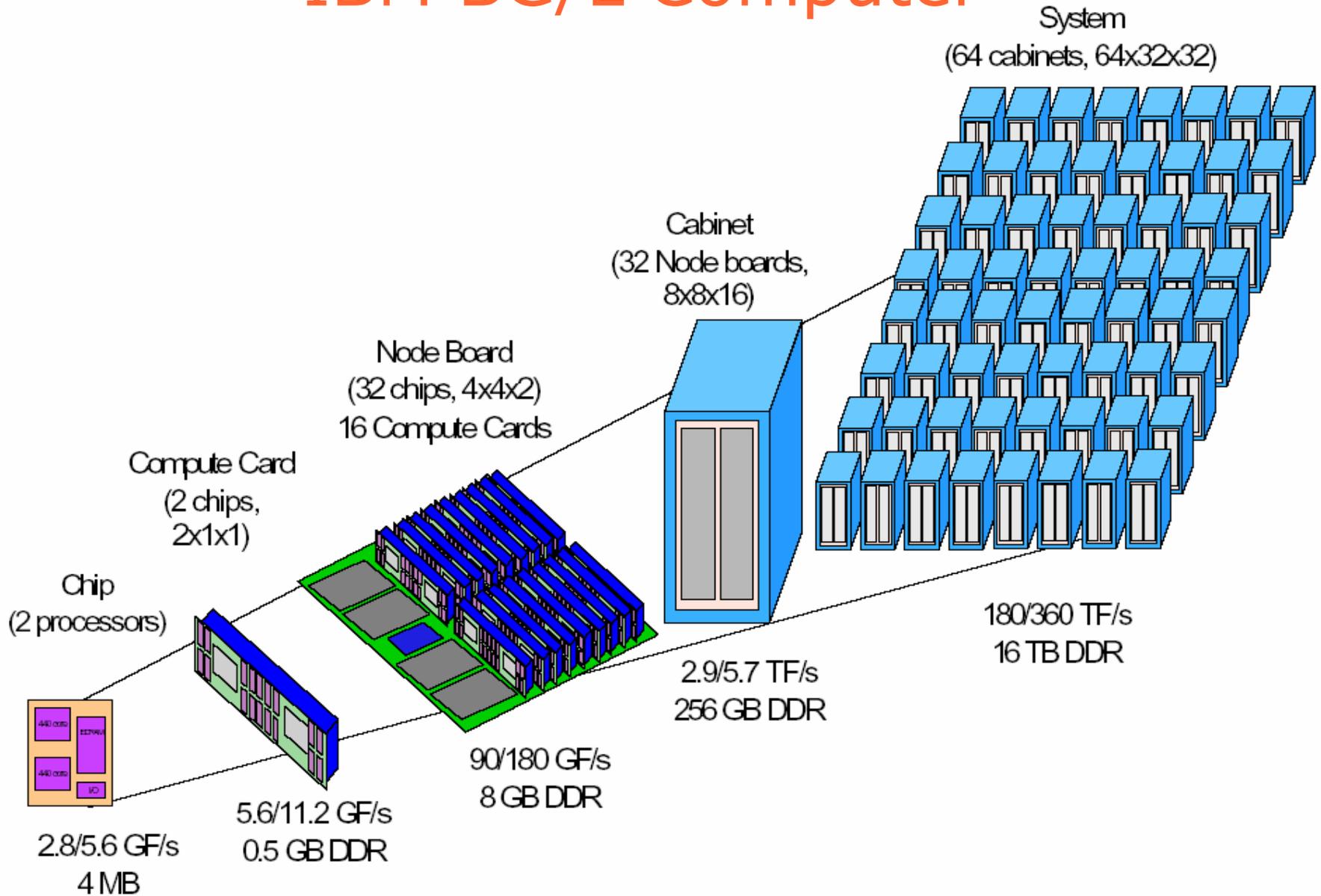


Type Ia Supernova Explosion:
Gravitationally Confined Detonation
(Calder, Plewa, Vladimirova, Lamb,
and Truran, 2004)





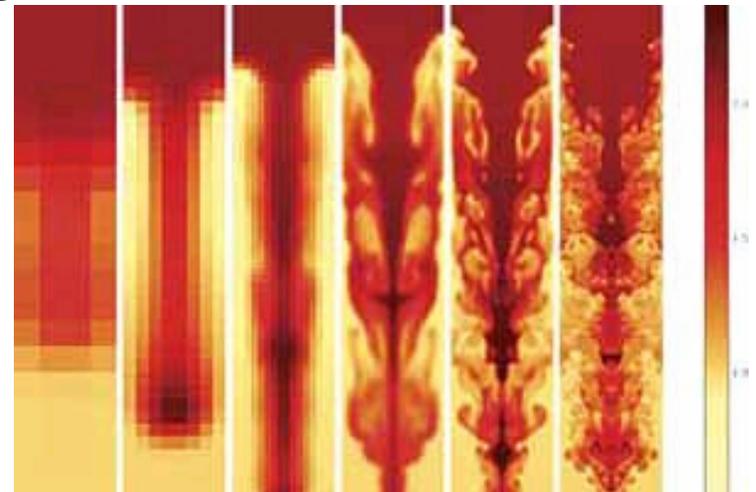
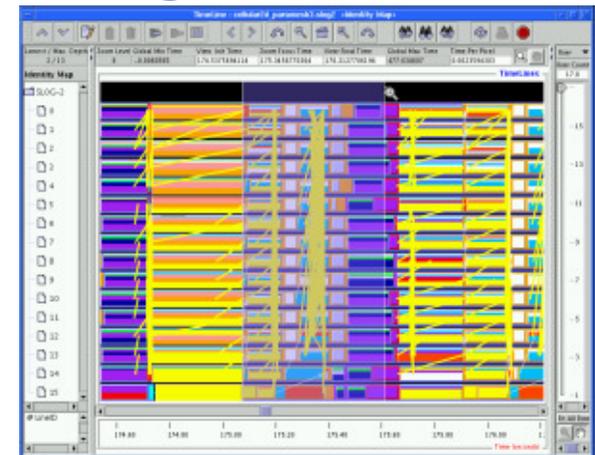
IBM BG/L Computer





Challenges Include ...

- Multi-scale, multi-physics modeling
 - ◆ Adaptive mesh refinement
 - ◆ Component architectures
- Scaling to 100K+ processors
 - ◆ Scalable parallel libraries
 - ◆ Parallel operating systems
- Understand & validating results
 - ◆ Visualization, data mining
 - ◆ Quantifying uncertainty





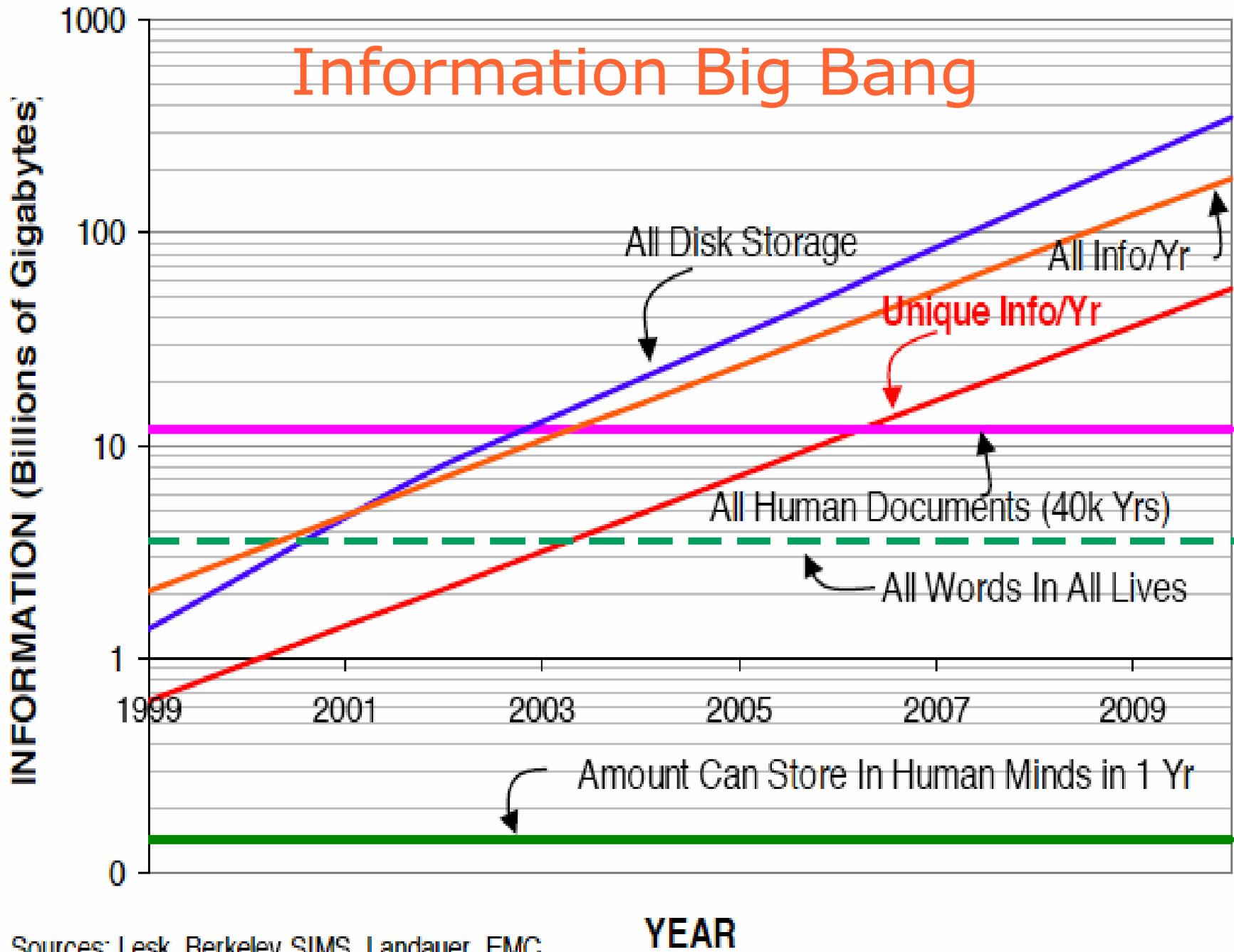
How Much Data?

- In 2006:
 - ◆ The world created 161 exabytes (1.6×10^{20} bytes) of digital data
 - ◆ There were one billion devices able to capture digital images
- By 2010:
 - ◆ Annual data output will reach one zettabyte (1×10^{21} bytes)



Source: IDC, 2007

Information Big Bang



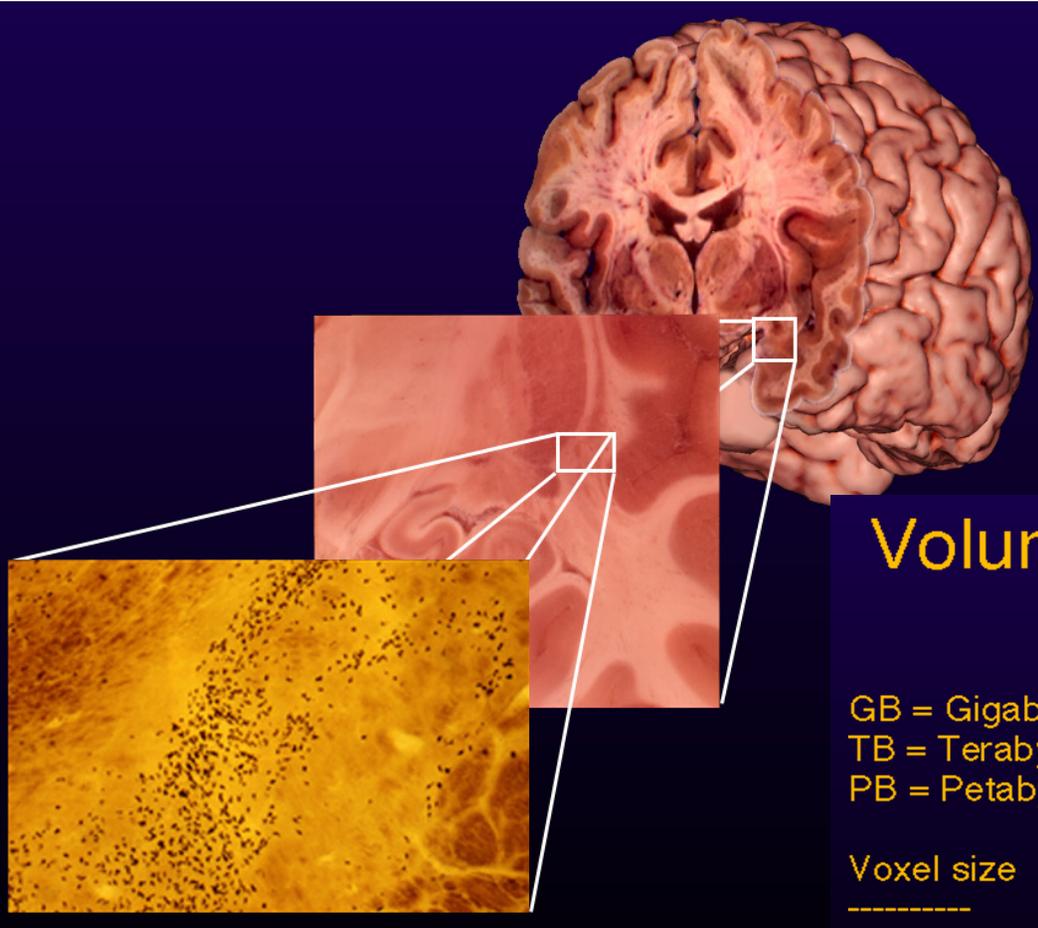
Sources: Lesk, Berkeley SIMS, Landauer, EMC

A Data Deluge



A Brain is a Lot of Data!

(Mark Ellisman, UCSD)



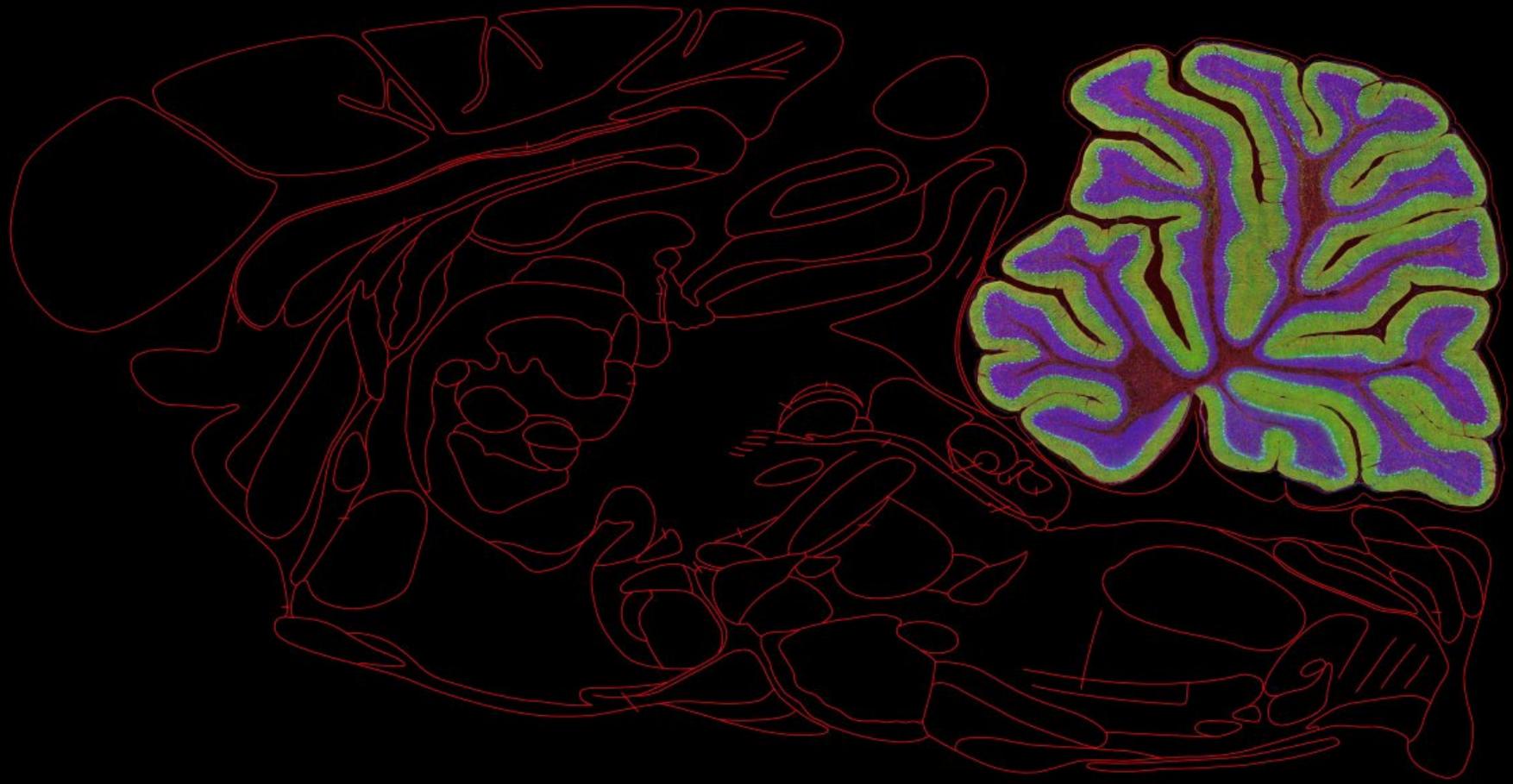
Volume sizes by resolution - brain = 1500 cm³

GB = Gigabyte = 10⁹
TB = Terabyte = 10¹²
PB = Petabyte = 10¹⁵

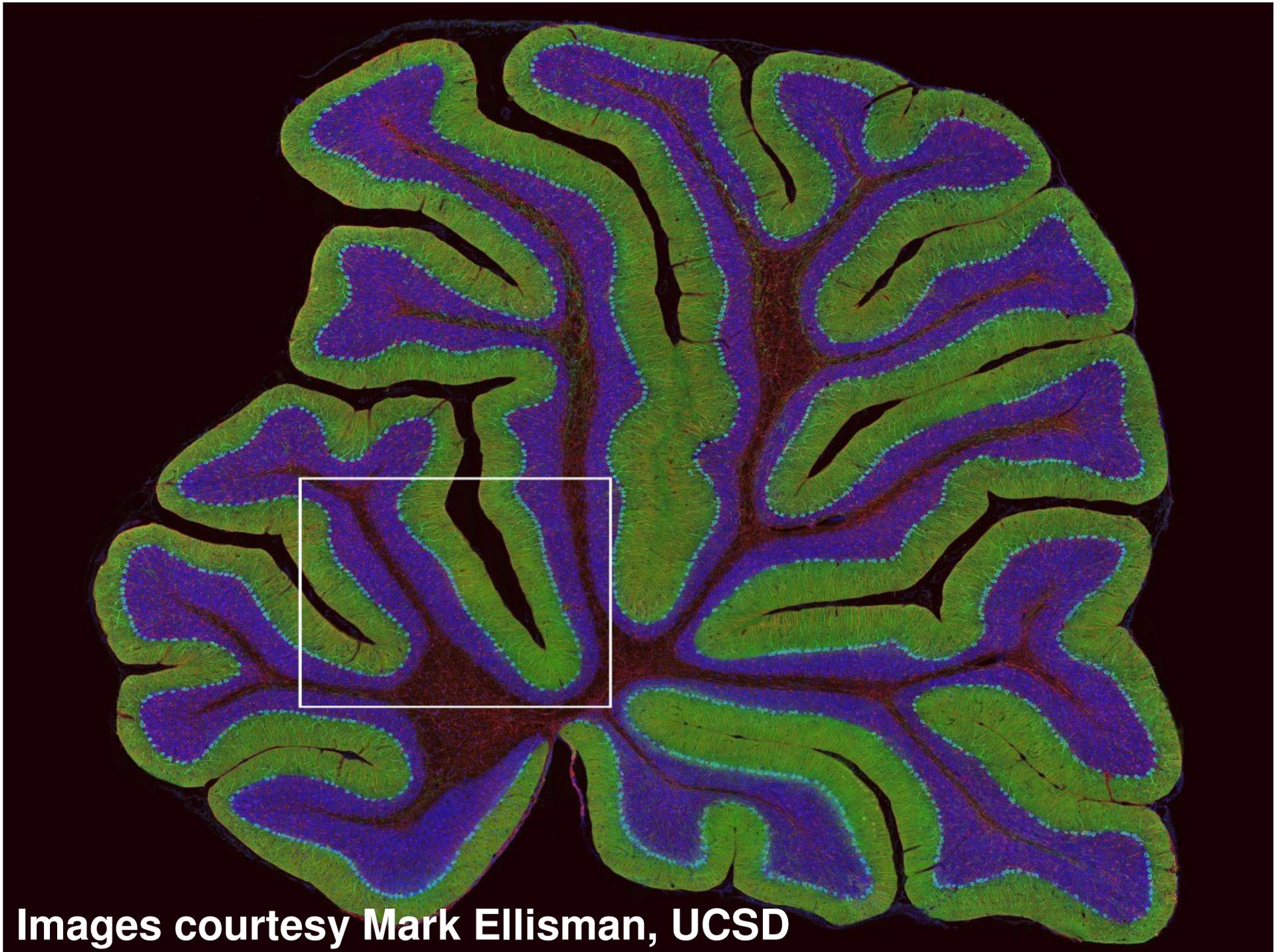
Voxel size	B&W (1 B/p)	High res (2 B/p)	Color (3 B/p)
cm	1.5 KB	3 KB	4.5 KB
mm	1.5 MB	3 MB	4.5 MB
10 μ m	1.5 TB	3 TB	4.5 TB
μ m	1.5 PB	3 PB	4.5 PB

*And comparisons must be
made among many*

We need to get to one micron to know location of every cell. We are starting to get to 10 microns



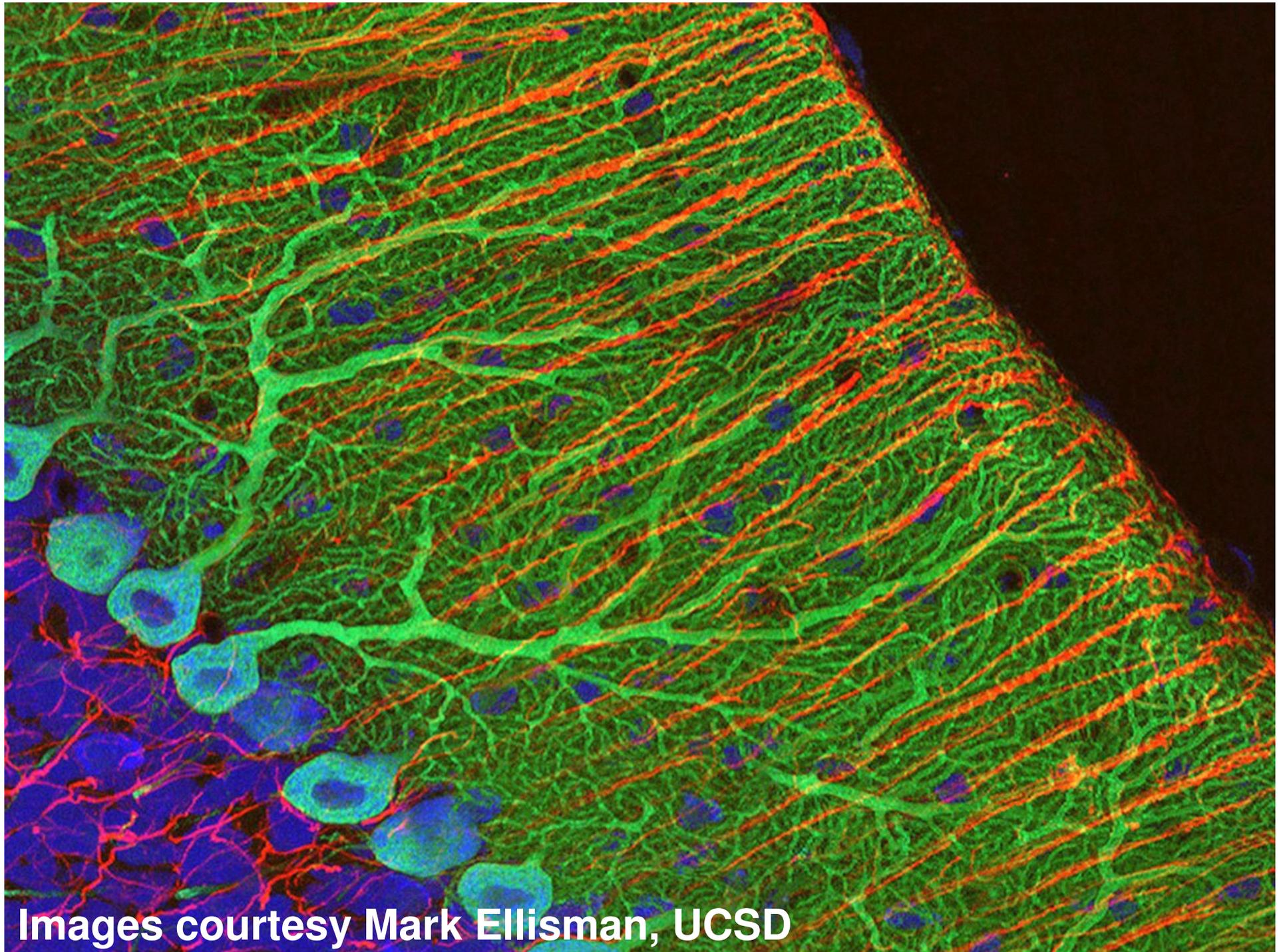
Images courtesy Mark Ellisman, UCSD



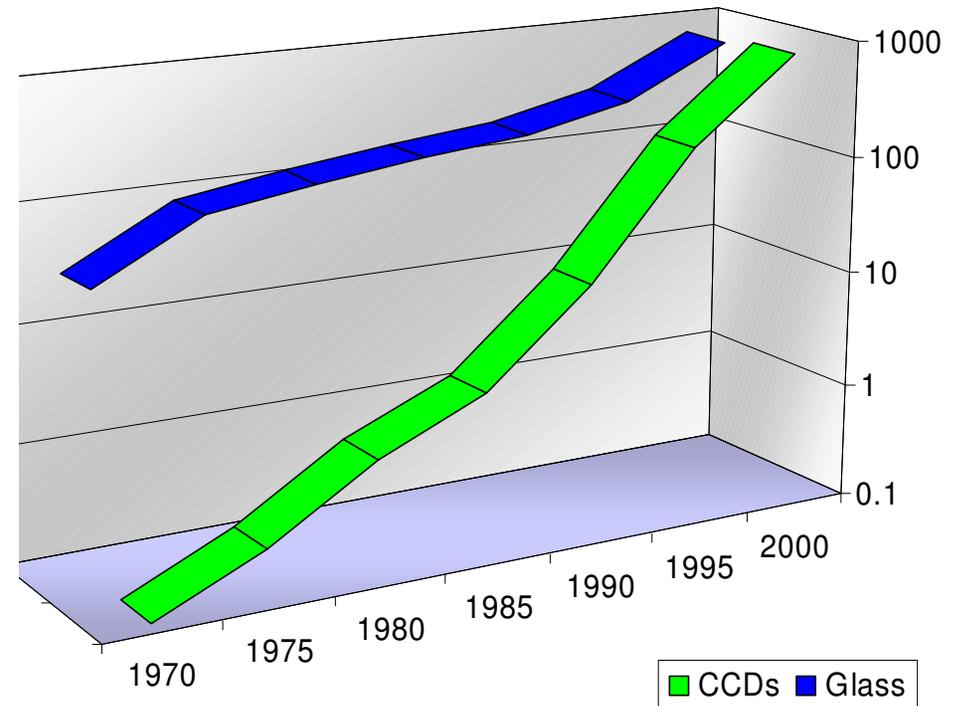
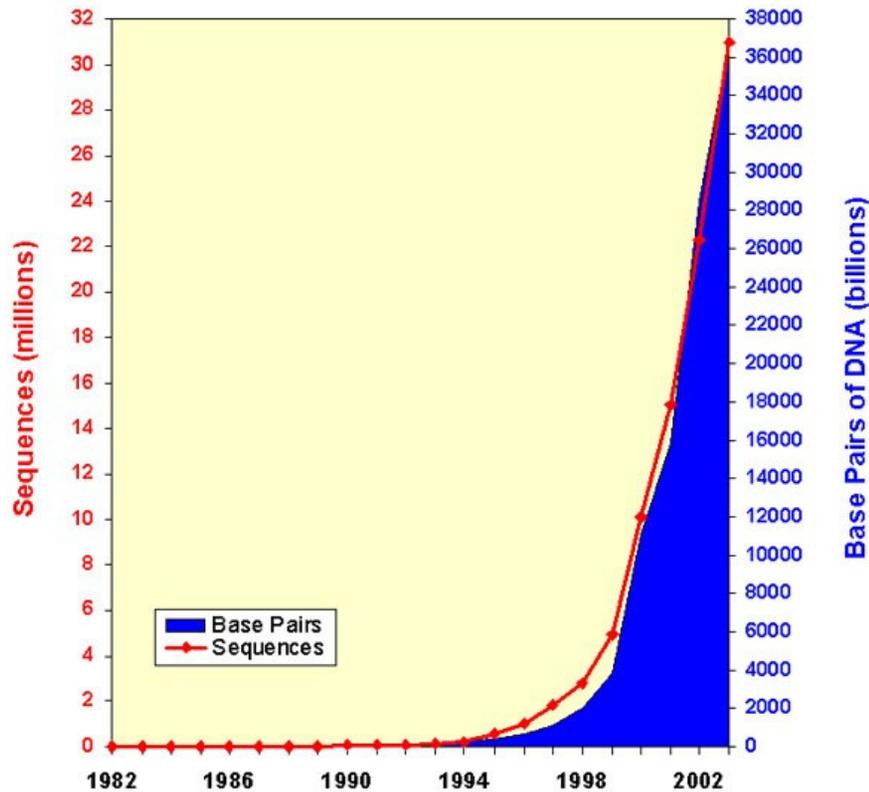
Images courtesy Mark Ellisman, UCSD



Images courtesy Mark Ellisman, UCSD



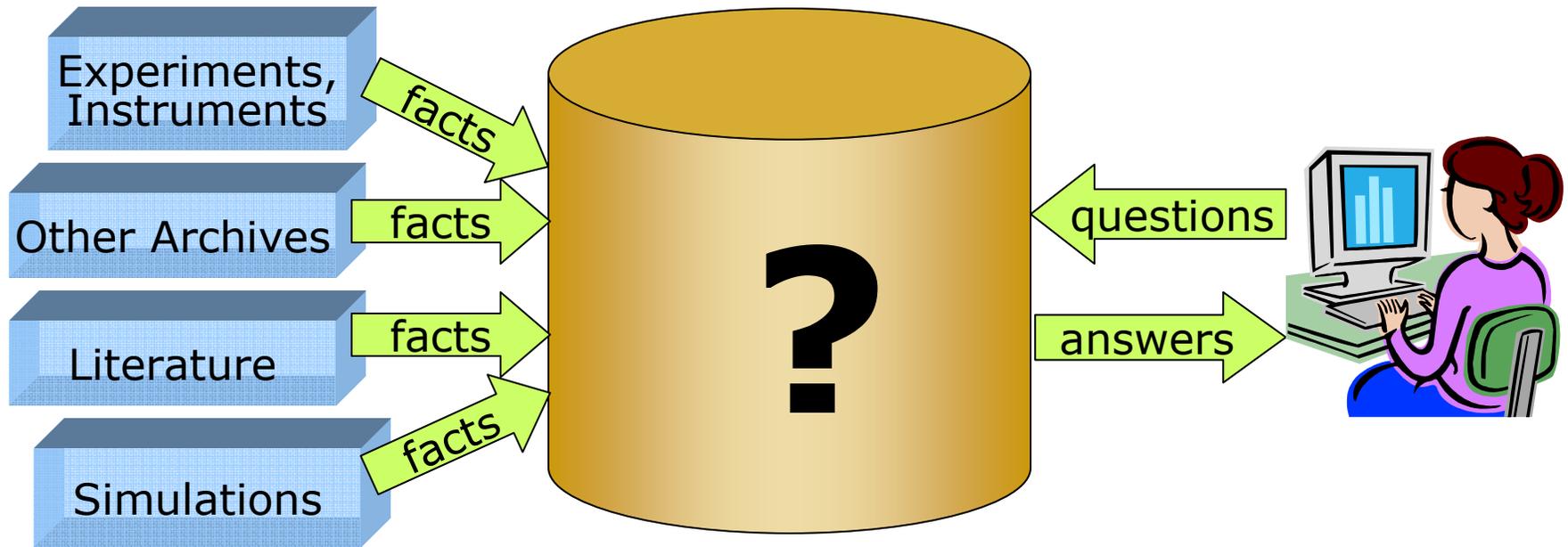
Images courtesy Mark Ellisman, UCSD



- Understanding increases **far** more slowly
- Methodological bottlenecks?
 - ➔ Improved technology
- Human limitations?
 - ➔ AI-assisted discovery



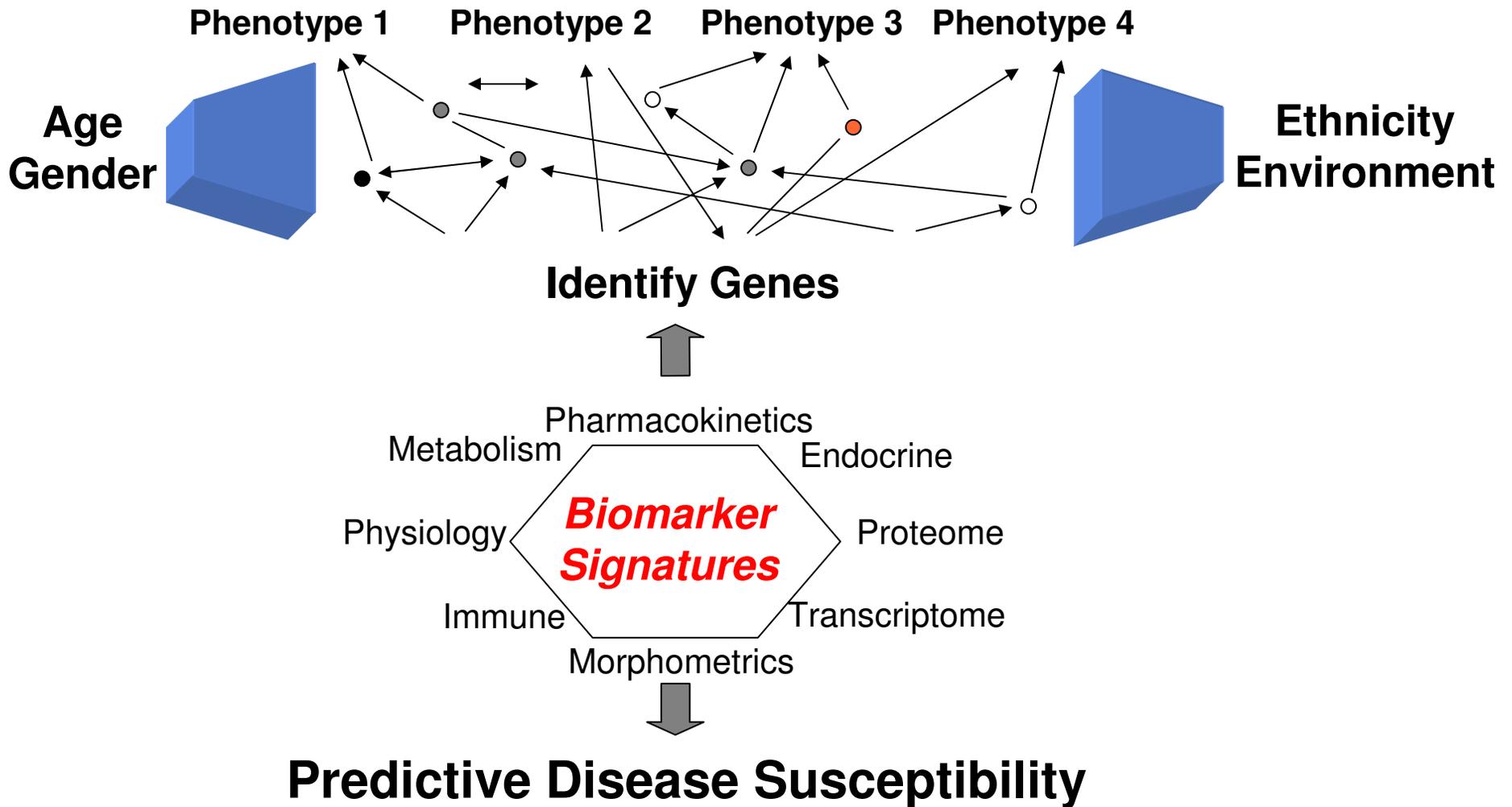
The Problem



- Data ingest
- Managing a petabyte
- Common schema
- How to organize it?
- How to *reorganize* it?
- Data query & visualization
- Support/training
- Performance: interactivity, scale in data size, analysis complexity, demand



Evidence Integration: Genetics & Disease Susceptibility





GeneWays as an Info-Grinder

On-line Journals

Andrey Rzhetsky et al.,
U.Chicago

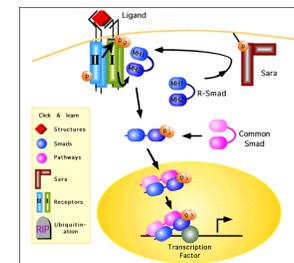
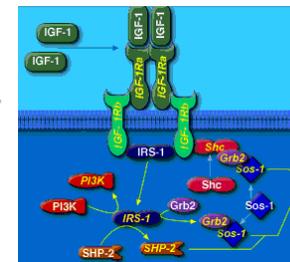
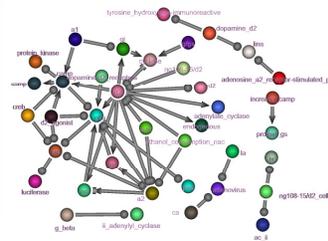


**We are currently screening
250,000 journal articles...**

2.5M reasoning chains

4M statements

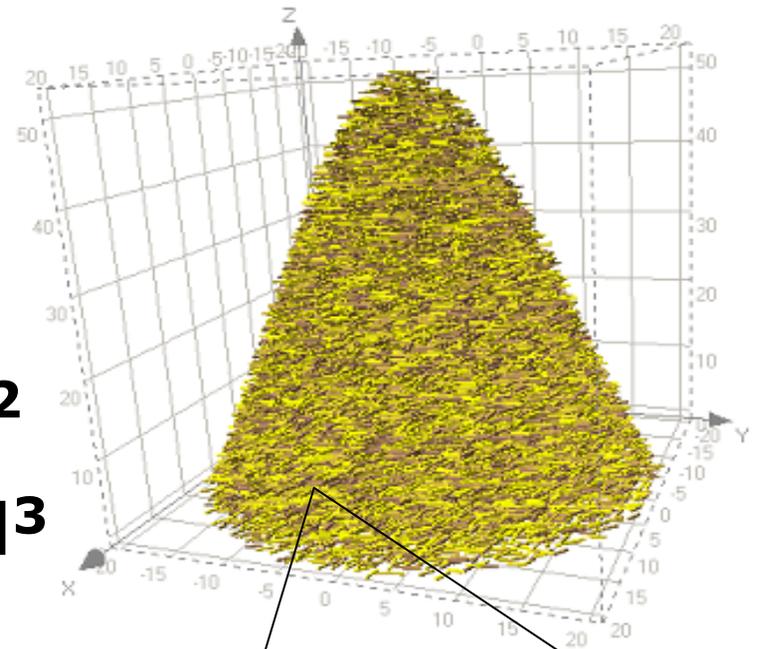
Pathways





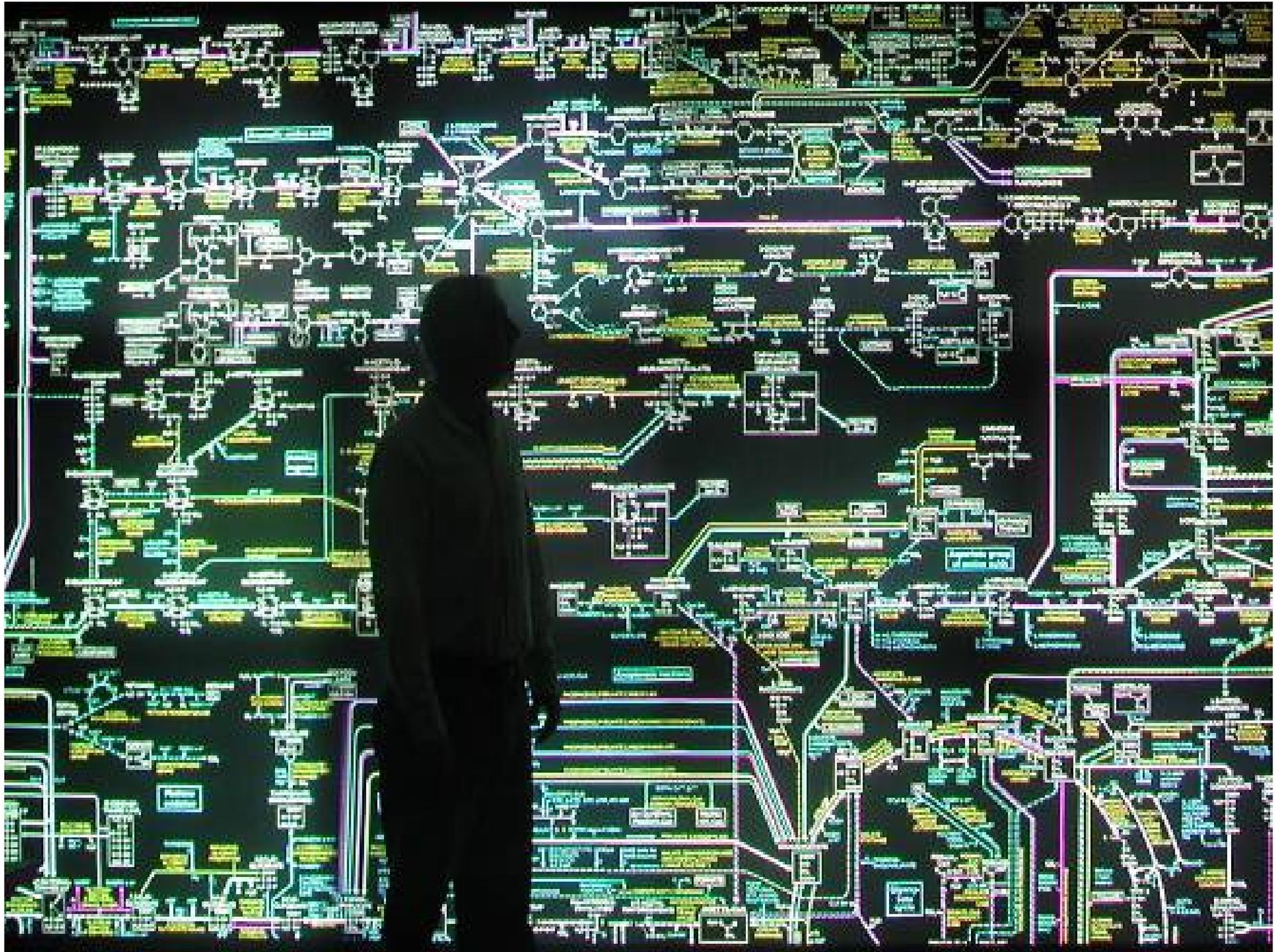
Data Analysis gets Fuzzy

- Global statistics?
 - ◆ Correlation functions: $\mathbf{N^2}$
 - ◆ Likelihood techniques: $\mathbf{N^3}$
- Best we can do is \mathbf{N} or maybe $\mathbf{N \log N}$



Haystack: Jim Gray/Alex Szalay

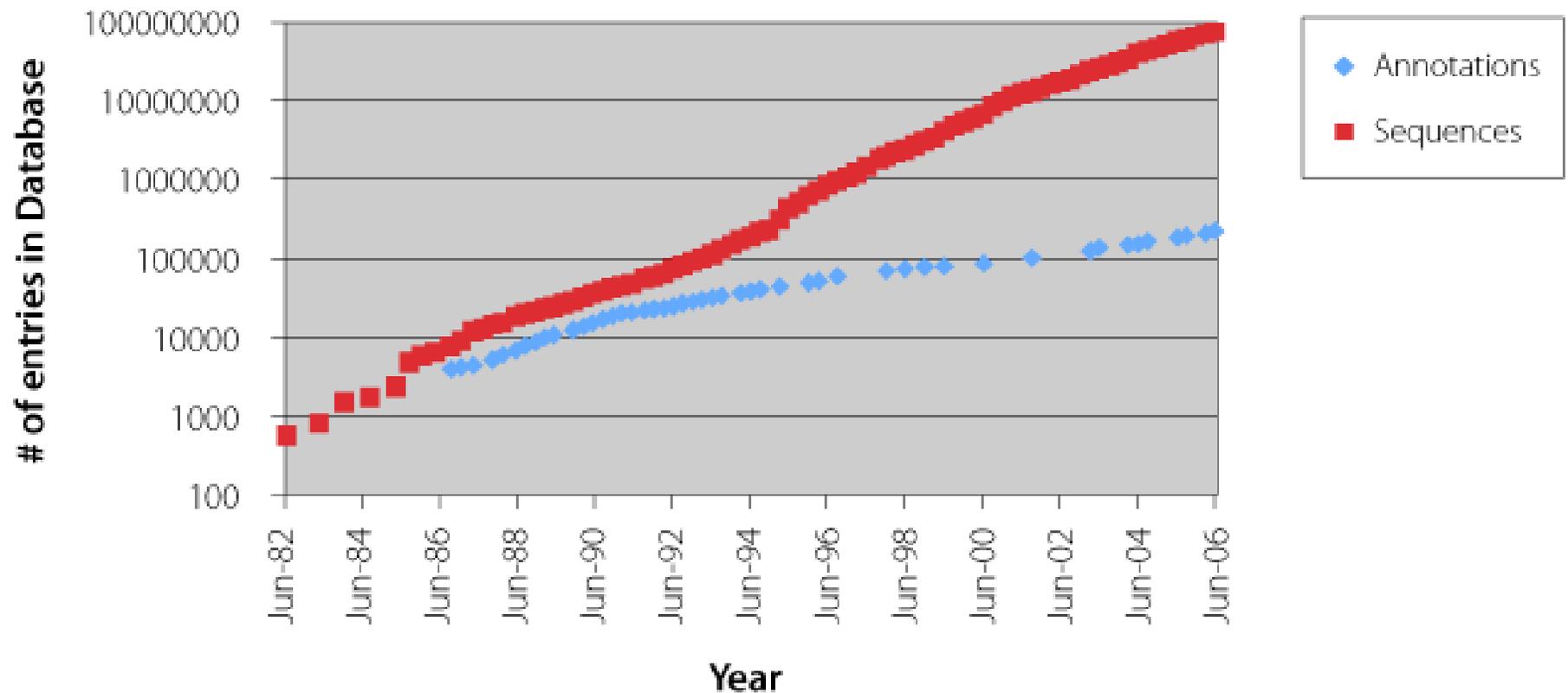
(scale approximate)²³





Growth of Sequences and Annotations since 1982

Growth of sequences and annotations since 1982

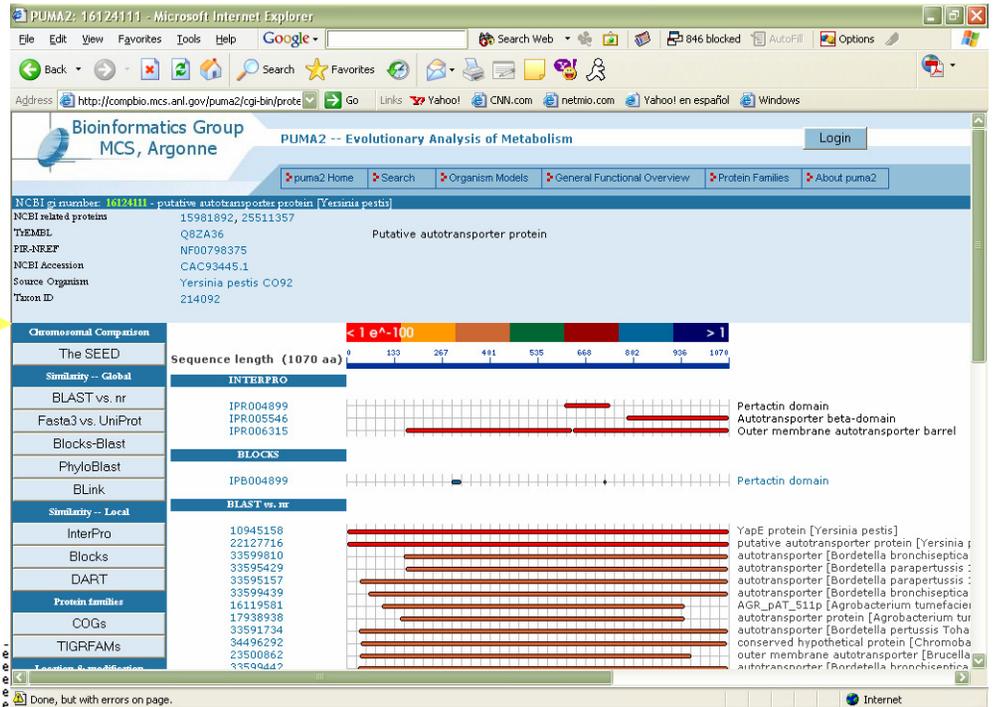


Folker Meyer, Genome Sequencing vs. Moore's Law: Cyber Challenges for the Next Decade, **CTWatch**, August 2006.



Production Science: Biology

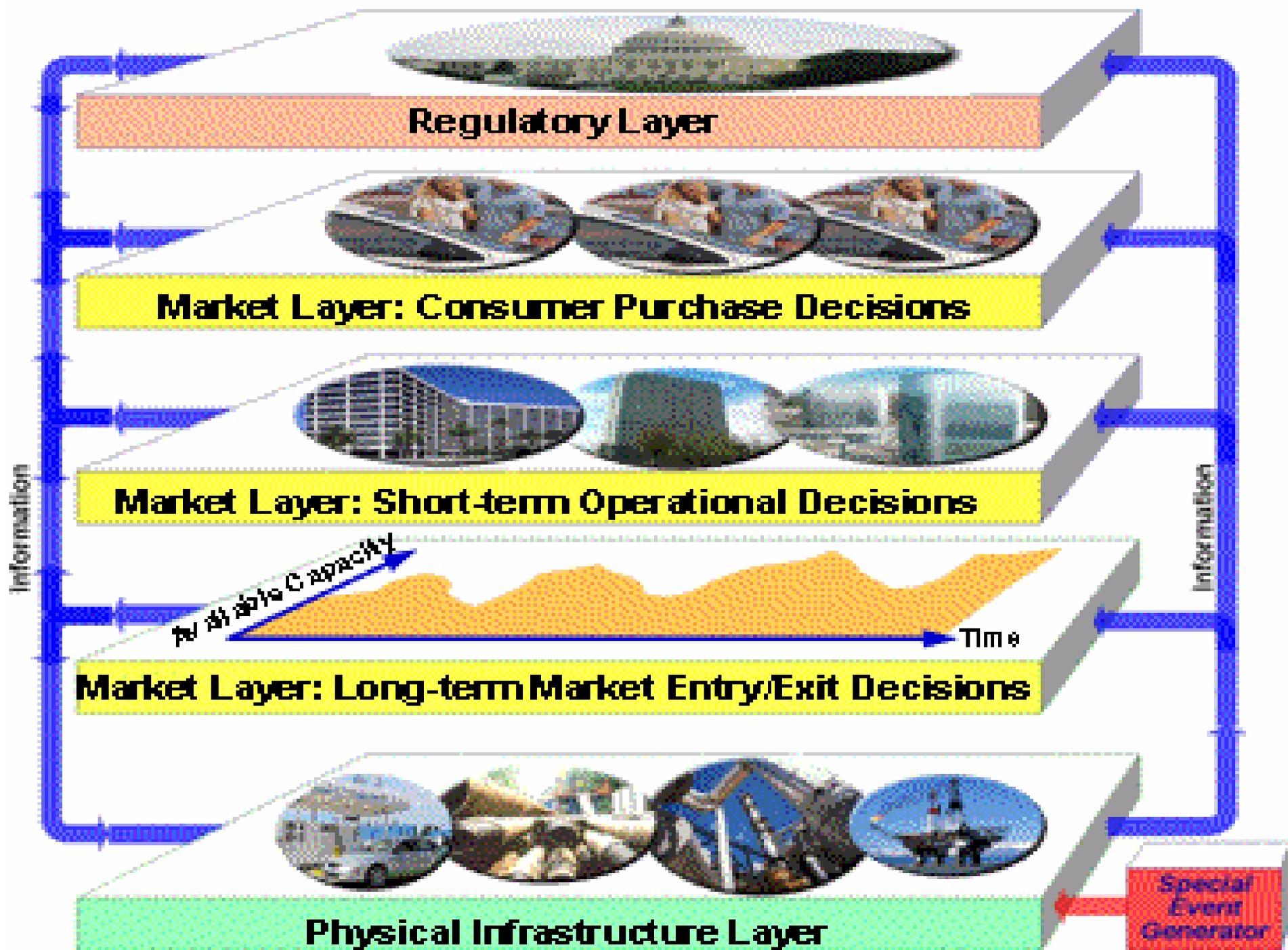
Public PUMA Knowledge Base
Information about proteins analyzed against ~2 million gene sequences



gi 23499780 gn REF_tigr BRA0013	gi 16080253 ref NP_391080.1	44.27	253	131	1	15	257	8	2603.7	e-30	134.4
gi 23499780 gn REF_tigr BRA0013	gi 23098409 ref NP_691875.1	43.48	253	133	2	16	258	5	2573.8	e-30	134.4
gi 23499780 gn REF_tigr BRA0013	gi 48637187 ref ZP_00294182.1	44.92	256	125	2	14	256	7	2591.1	e-30	134.4
gi 23499780 gn REF_tigr BRA0013	gi 52005400 gb AA025342.1	44.75	257	126	2	15	258	3	2561.9	e-30	134.4
gi 23499780 gn REF_tigr BRA0013	gi 48664015 ref ZP_00317908.1	44.49	245	134	1	13	257	5	2476.1	e-30	134.4
gi 23499780 gn REF_tigr BRA0013	gi 30348891 gb AA028934.1	39.53	253	138	3	18	257	5	2522.0	e-43	177.6
gi 23499780 gn REF_tigr BRA0013	gi 19655222 gb AA093939.1	40.64	251	138	1	17	256	10	2602.7	e-43	177.6
gi 23499780 gn REF_tigr BRA0013	gi 27358806 gb AA007757.1	43.03	251	130	4	18	256	11	2602.5	e-41	170.6
gi 23499780 gn REF_tigr BRA0013	gi 12597924 gb AA018599.2	46.70	162	96	1	62	243	5	1856.8	e-39	162.5
gi 23499780 gn REF_tigr BRA0013	gi 46363318 ref ZP_0026079.1	39.58	240	135	2	14	253	6	2361.8	e-36	154.5
REF_tigr BRA0013	gi 39933731 ref NP_946007.1	34.90	255						e-33	142.9	
REF_tigr BRA0013	gi 48782600 ref ZP_00279106.1	35.92	245						e-32	141.4	
REF_tigr BRA0013	gi 41407534 ref NP_960370.1	36.09	266						e-32	139.4	
REF_tigr BRA0013	gi 48851585 ref ZP_00305793.1	32.39	247						e-32	139.0	
REF_tigr BRA0013	gi 15966306 ref NP_386659.1	36.50	263						e-31	137.9	
REF_tigr BRA0013	gi 17548526 ref NP_521866.1	36.36	264						e-31	137.1	
gi 23499780 gn REF_tigr BRA0013	gi 51891730 ref VP_074421.1	38.87	247	136	7	18	256	1	2403.4	e-30	133.7
gi 23499780 gn REF_tigr BRA0013	gi 145881 gb AA23739.1	33.87	246	147	3	13	253	3	2404.4	e-30	133.3
gi 23499780 gn REF_tigr BRA0013	gi 25029334 ref NP_739388.1	35.20	250	147	4	15	256	6	2485.7	e-30	132.9
gi 23499780 gn REF_tigr BRA0013	gi 21220953 ref NP_536732.1	36.52	257	138	6	12	255	5	2545.7	e-30	132.9
gi 23499780 gn REF_tigr BRA0013	gi 46314029 ref ZP_00214635.1	33.86	254	153	2	12	258	3	2485.7	e-30	132.9
gi 23499780 gn REF_tigr BRA0013	gi 41406852 ref NP_959683.1	35.61	238	149	2	16	253	2	2309.8	e-30	132.1
gi 23499780 gn REF_tigr BRA0013	gi 15564471 ref NP_229523.1	35.69	255	144	5	12	256	2	2469.8	e-30	132.1
gi 23499780 gn REF_tigr BRA0013	gi 23470090 ref ZP_00125423.1	35.20	250	145	4	12	253	3	2439.8	e-30	132.1
gi 23499780 gn REF_tigr BRA0013	gi 24935279 gb AA064257.1	34.63	257	146	4	12	257	4	2499.8	e-30	132.1
gi 23499780 gn REF_tigr BRA0013	gi 48347651 ref ZP_00301815.1	36.05	258	145	9	12	257	4	2531.3	e-29	131.7
gi 23499780 gn REF_tigr BRA0013	gi 28851510 gb AA054587.1	36.40	250	142	4	12	253	3	2431.3	e-29	131.7
gi 23499780 gn REF_tigr BRA0013	gi 12737873 ref NP_770312.1	36.25	251	143	3	14	255	7	2491.3	e-29	131.7
gi 23499780 gn REF_tigr BRA0013	gi 1708836 sp P50198 LIDX_PSEPA	34.23	260	143	4	12	257	4	2491.7	e-29	131.3
gi 23499780 gn REF_tigr BRA0013	gi 33594148 ref NP_381792.1	34.17	240	148	5	18	256	6	2363.7	e-29	130.2
gi 23499780 gn REF_tigr BRA0013	gi 33595116 ref NP_381759.1	34.17	240	148	5	18	256	6	2363.7	e-29	130.2
gi 23499780 gn REF_tigr BRA0013	gi 3328306 ref NP_232830.1	34.20	241	143	5	18	256	6	2363.7	e-29	130.2

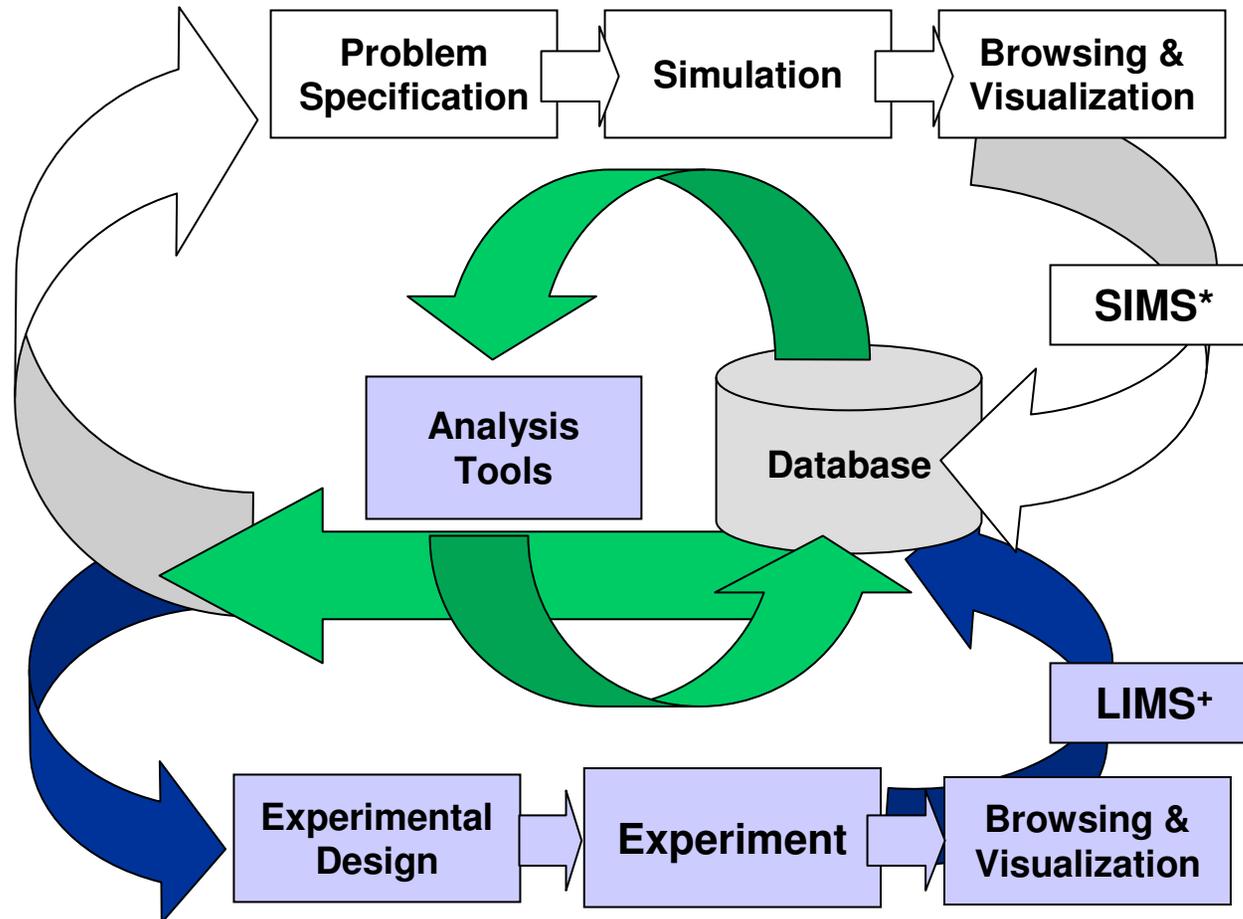
Back Office Analysis on Grid
Millions of BLAST, BLOCKS, etc., on OSG and TeraGrid

Natalia Maltsev et al., <http://compbio.mcs.anl.gov/puma2>





Integrated View of Simulation, Experiment, & Bioinformatics



*Simulation Information Management System

+Laboratory Information Management System



eScience

Computational science

+ Informatics

= eScience [John Taylor, UK EPSRC]

- ◆ “Large-scale science carried out through distributed collaborations—often leveraging access to large-scale data & computing”

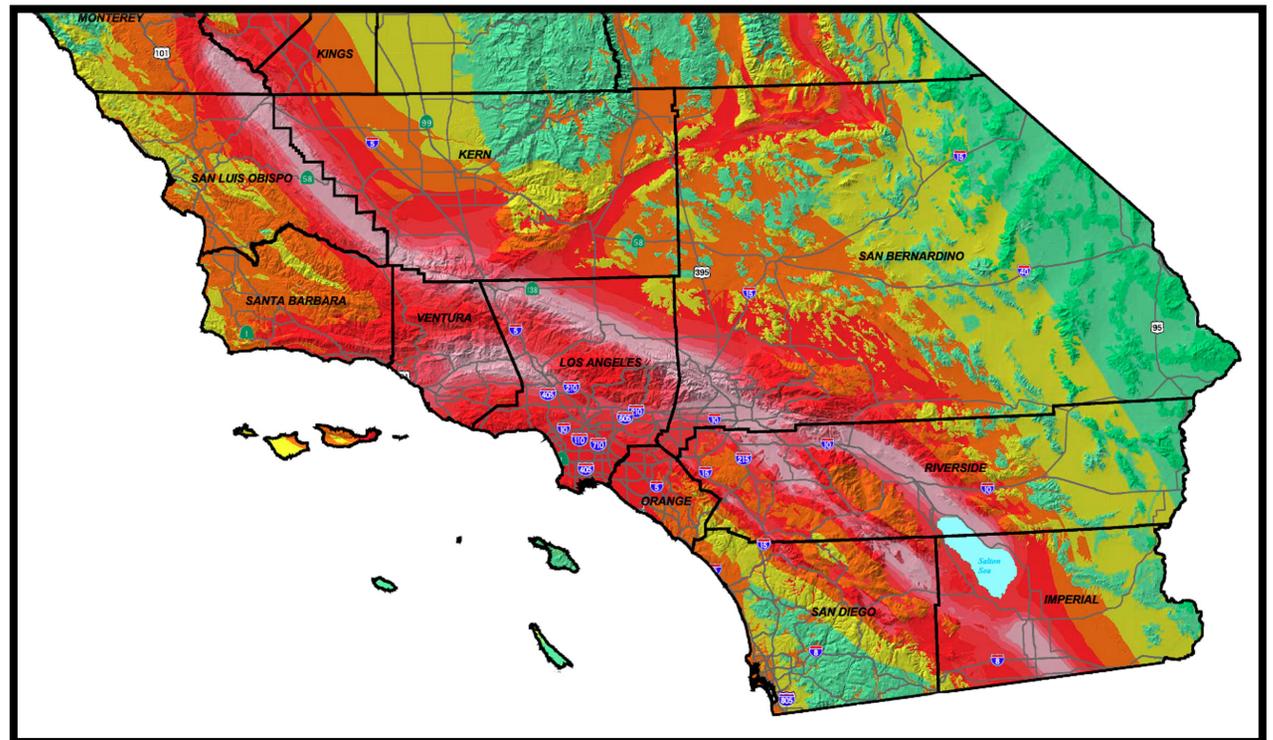


Seismic Hazard Analysis

Defn: Max. intensity of shaking expected at a site during a fixed time interval

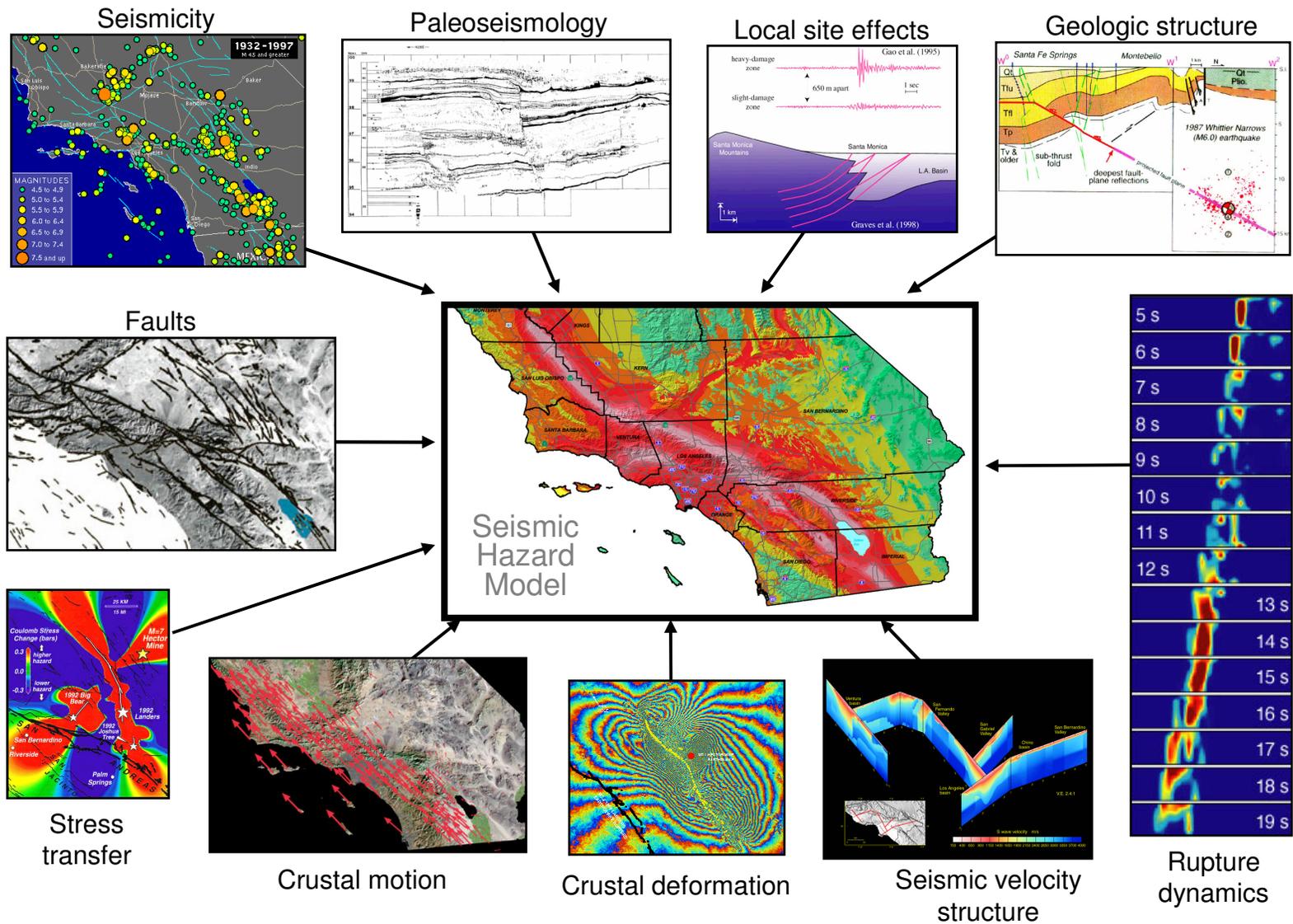
Example: National seismic hazard maps

- Intensity measure: peak ground acceleration
- Interval: 50 yrs
- Probability of exceedance: 2%





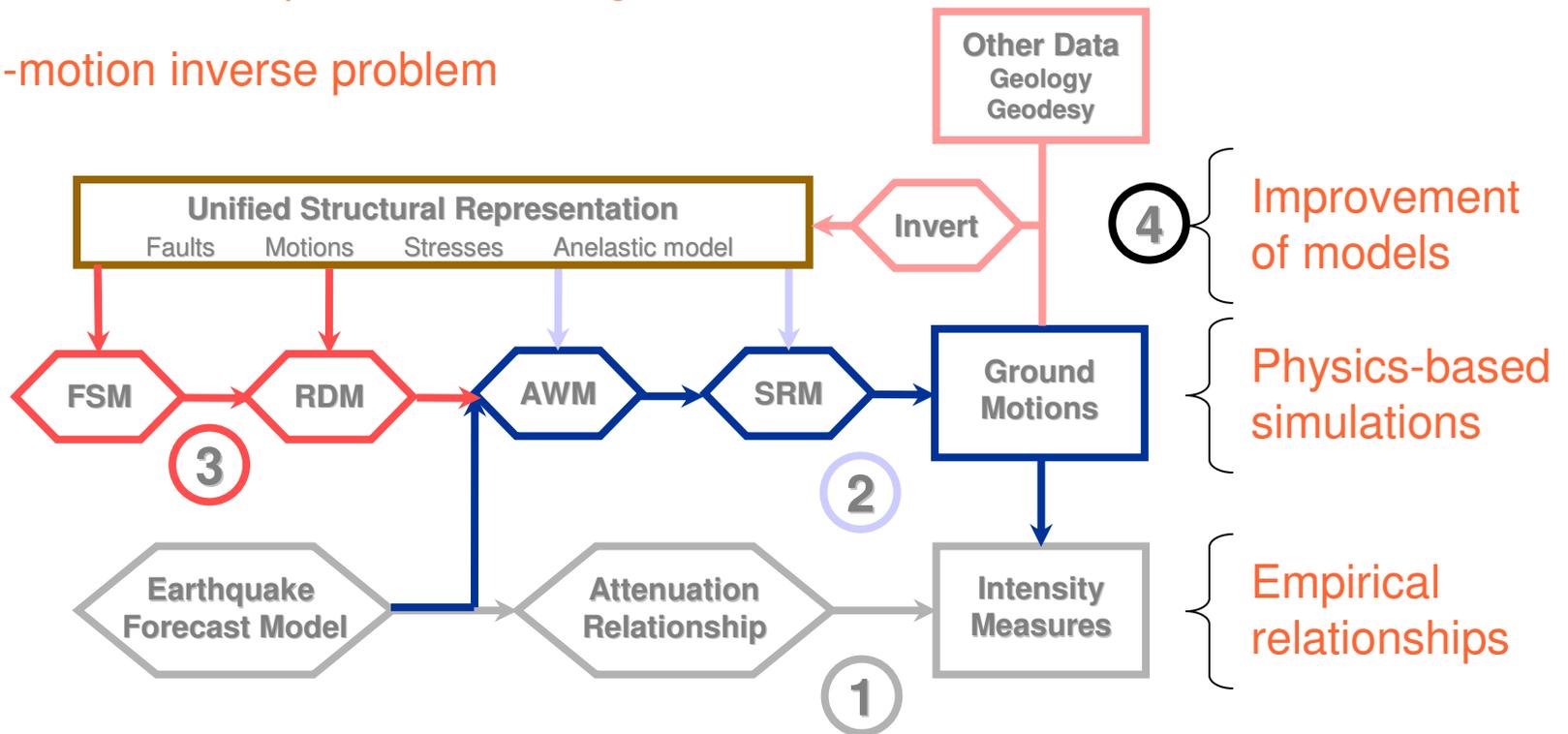
Seismic Hazard Analysis





- 1 Standardized Seismic Hazard Analysis
- 2 Ground motion simulation
- 3 Physics-based earthquake forecasting
- 4 Ground-motion inverse problem

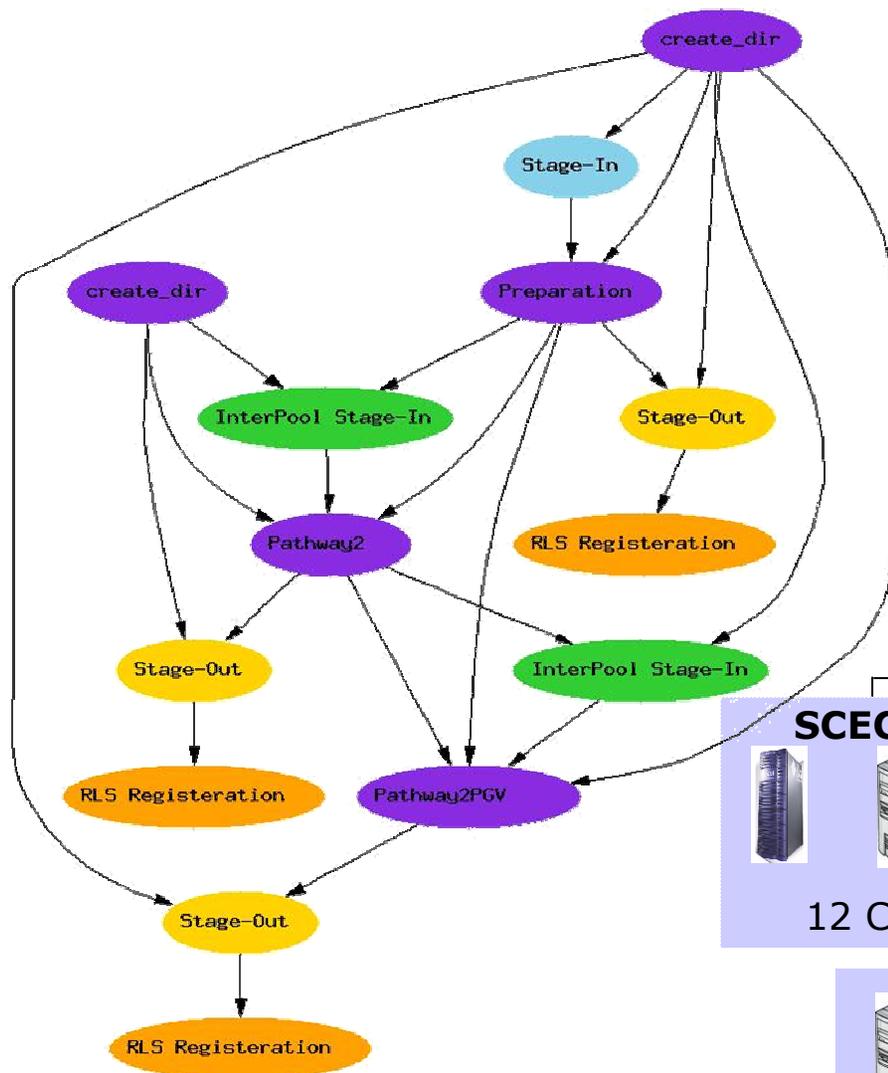
SHA Computational Pathways



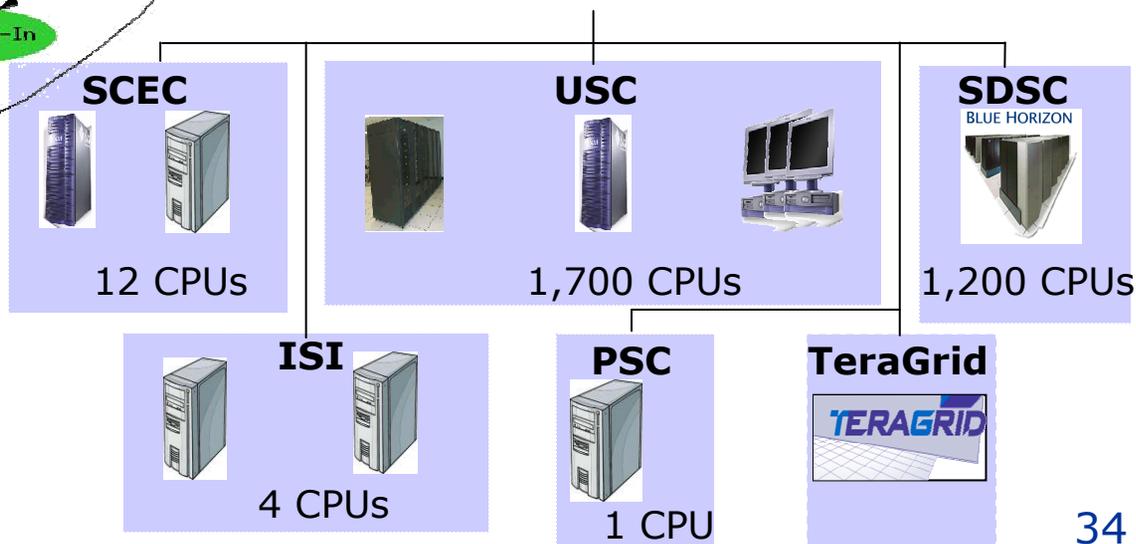
FSM = Fault System Model
RDM = Rupture Dynamics Model

AWP = Anelastic Wave Propagation
SRM = Site Response Model

Access to National Cyberinfrastructure



- Prepare input to Pathway2 wave propagation code
- Pathway2PGV converts output into hazard map
- Map is visualized



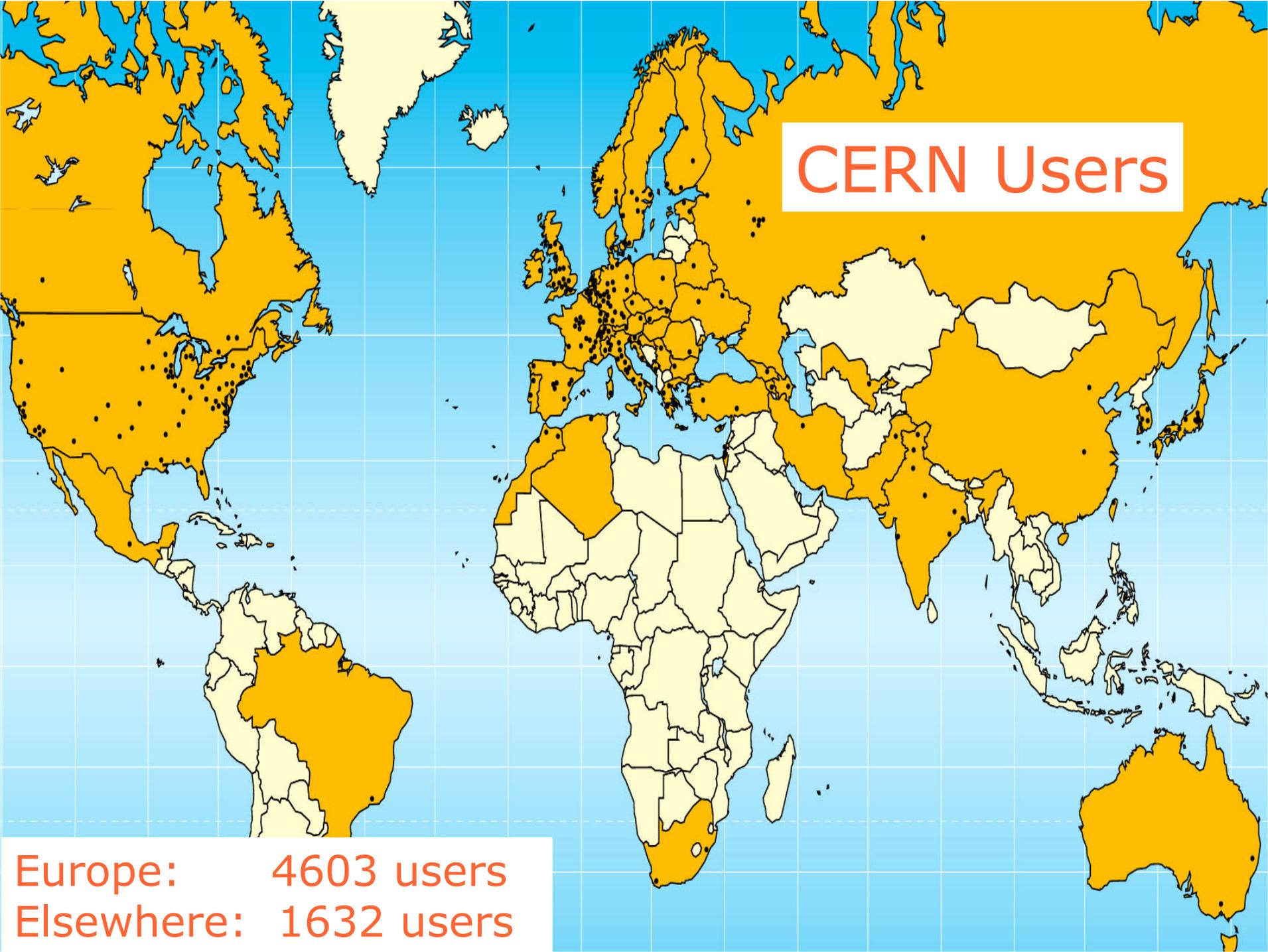
Slide 34

CFK11 This shows only pathway two, where the other pathways involved as well.

Do you have a visualization of the output of wone of these runs?

Are the CPU pictures accurate to what you ran?

Carl Kesselman, 9/27/2004

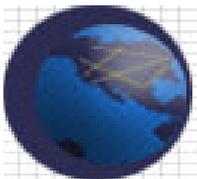
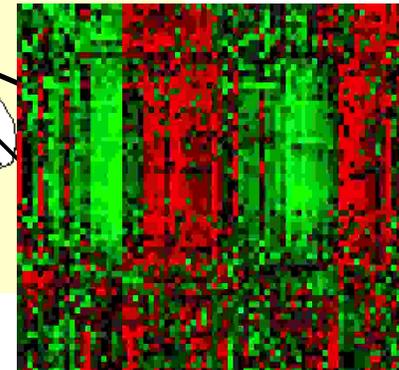
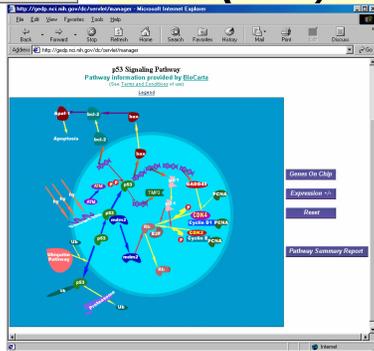
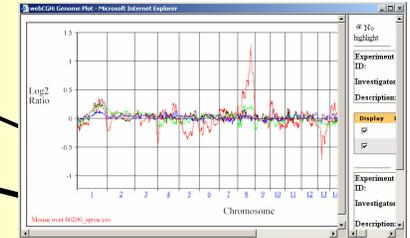
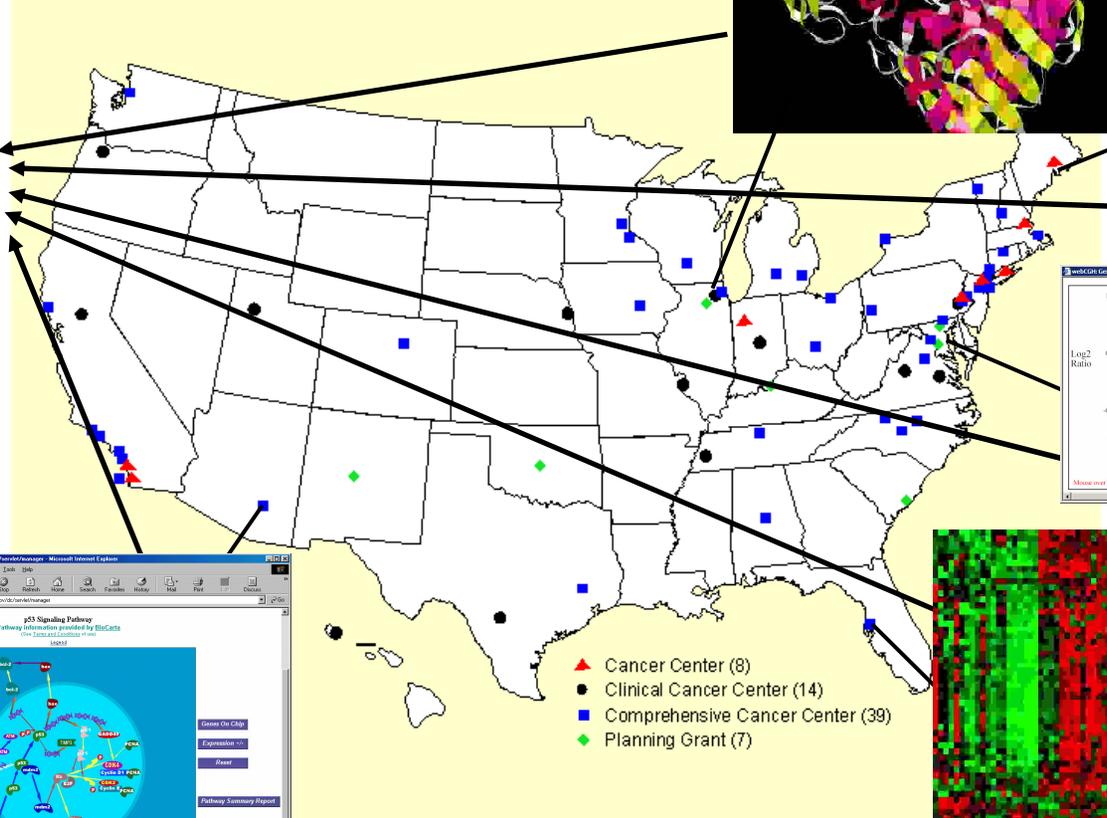
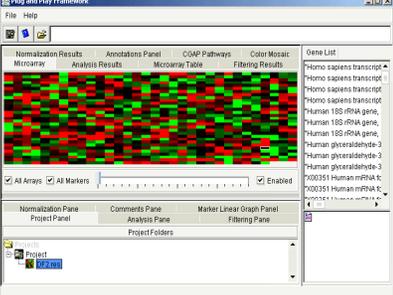
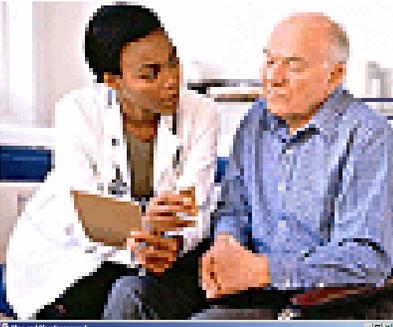
A world map with a light blue grid background. Landmasses are colored in shades of orange and yellow. Small black dots representing CERN users are scattered across the map, with a high concentration in Europe and North America. A white rectangular box with a red border is positioned in the upper right quadrant of the map, containing the text 'CERN Users'.

CERN Users

Europe: 4603 users
Elsewhere: 1632 users



Cancer Biology

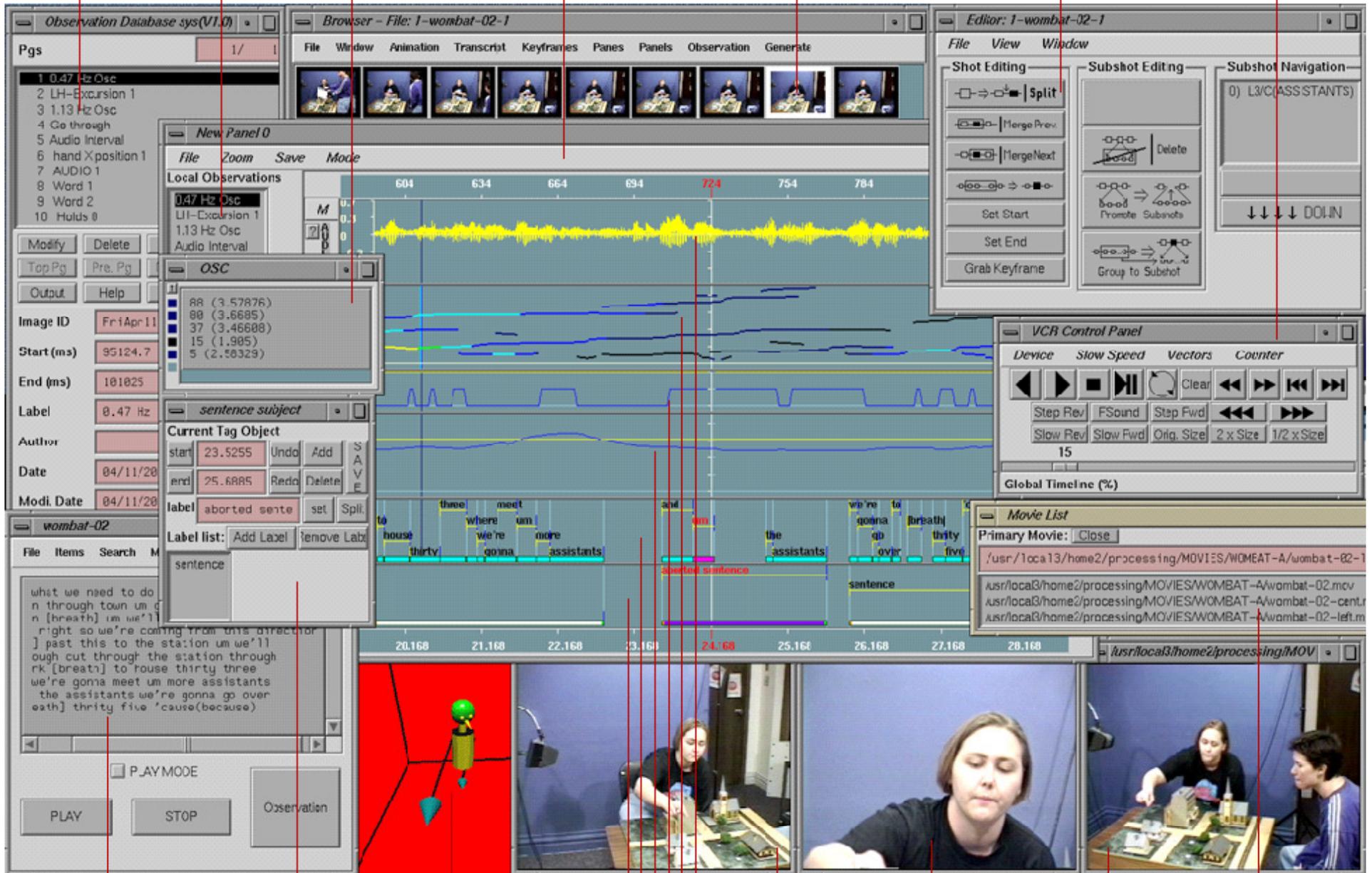


caBIG

cancer Biomedical
Informatics Grid



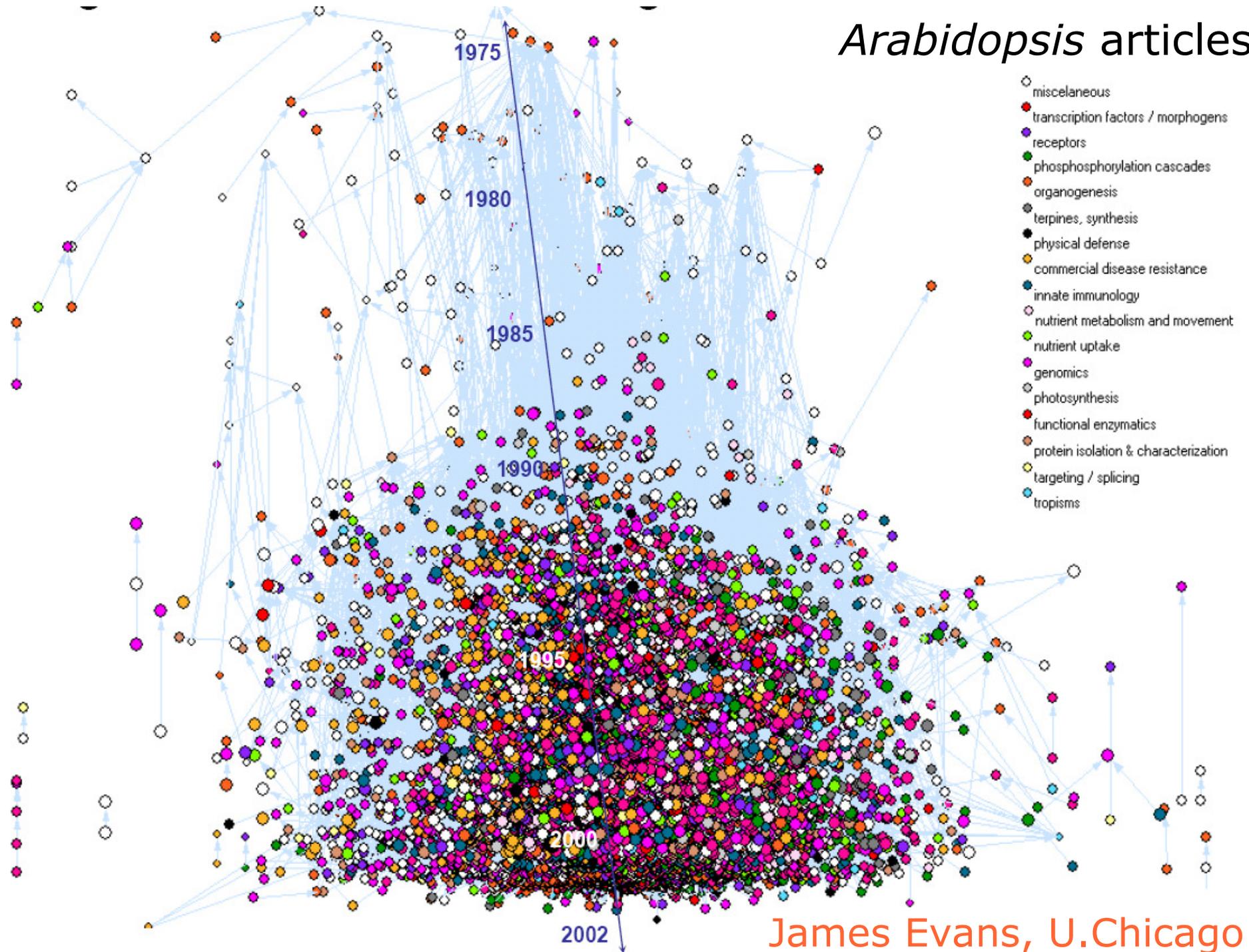
Global Observation Database (View) Graph-Based Observations Information Window Animated Graph Representation (Graph Panel) Hierarchical Shot-Keyframe Representation Shot-Keyframe Hierarchy Editor VCR-Style Control Panel



Animated Text Transcript (Paragraph Representation) Tag Transcript Editor Animated Avatar Representation Animated Graph Panes Video Displays Video List

Bennett Berthenthal et al., www.sidgrid.org

Arabidopsis articles



James Evans, U.Chicago



eScience Challenges

- Simulate complex, multi-component systems
- Evaluate accuracy of such simulations
- Integrate evidence to draw conclusions
- Evaluate strength of conclusions
- Automate “experimental” workflows
- Document basis for conclusions (provenance)
- Allow these problems to be tackled by distributed teams using federated resources



What is Fundamental?

Bits

+

Algorithms

+

Complex systems

First two, at least © CS

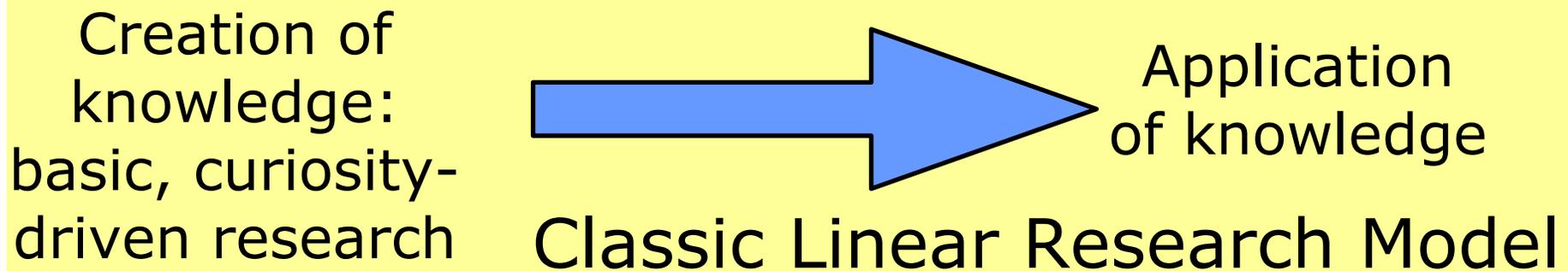
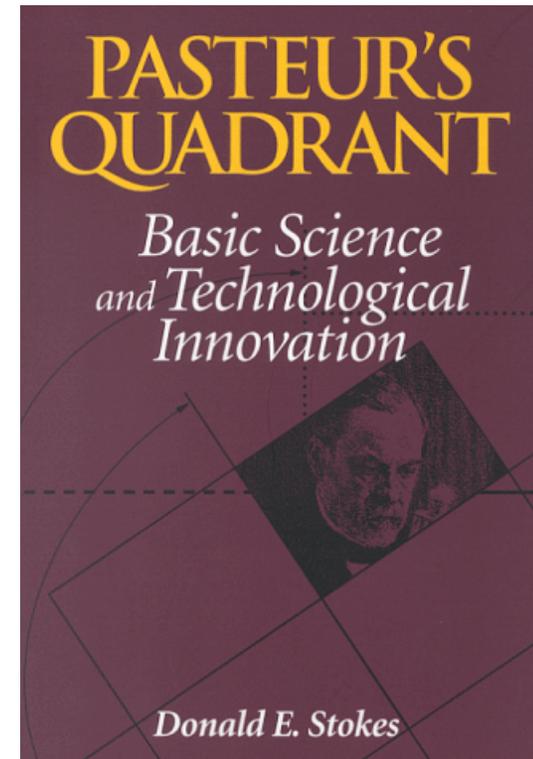
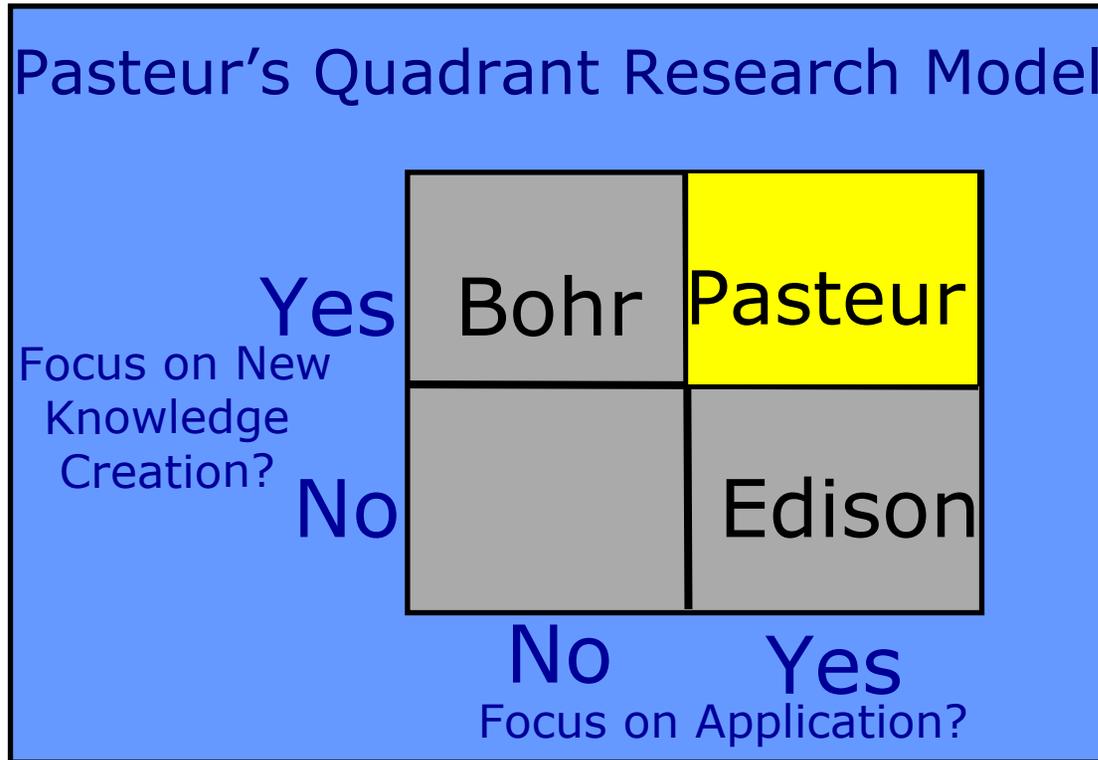


Computer Science: A Narrow or Broad View?

- **Narrow**
 - ◆ CS is programming
 - ➔ Other aspects of information are the domain of “statistics,” “bioinformatics”, etc., etc.
- **Broad**
 - ◆ CS is “the systematic study of algorithmic processes that describe and transform information, their theory, analysis, design, efficiency, implementation, and application” (Denning et al., CACM, 1989)
 - ➔ Statistics & bioinformatics are subdisciplines of computer science



Effective eScience **requires** PQ research models



Slide courtesy Dan Atkins, U.Michigan



“Applied computer science is now playing the role that mathematics did from the 17th through the 20th centuries: providing an orderly, formal framework & exploratory apparatus for other sciences.”

—George Djorgovski

“... the branch of computer science that concerns itself with the application of computing knowledge to other domains”?



Computation Institute

A joint institute of Argonne and the University of Chicago, focused on advancing **system-level science**

Solutions to many grand challenges facing science and society today are dependent upon the analysis and understanding of entire systems, not just individual components. They require not reductionist approaches but the synthesis of knowledge from multiple levels of a system, whether biological, physical, or social (or all three).

<http://www.ci.uchicago.edu>



Thanks!

- foster@mcs.anl.gov
- <http://www.ci.uchicago.edu>
- <http://ianfoster.typepad.com>



In Memoriam: Jim Gray (1944-2007?)

Turing Award, 1998

“for seminal contributions to database & transaction processing and technical leadership in system implementation”

